

## Supplementary Material for “A Model-averaging method for high-dimensional regression with missing responses at random”

Jinhan Xie<sup>a</sup> Xiaodong Yan<sup>b</sup> and Niansheng Tang<sup>a</sup>

*a*Yunnan Key Laboratory of Statistical Modeling and Data Analysis, Yunnan University

*b*School of Economics, Shandong University

### Supplementary Material

#### S1. Properties of penalized likelihood estimator $\hat{\gamma}$

In this section, we investigate the consistency, and oracle property of the penalized likelihood estimator  $\hat{\gamma}$  of  $\gamma$ . Without loss of generality, we can write  $\gamma = (\gamma_1^\top, \gamma_2^\top)^\top$ , where  $\gamma_1 \in \mathbf{R}^{d_m}$  and  $\gamma_2 \in \mathbf{R}^{q-d_m}$  correspond to the nonzero and zero components of  $\gamma$ , respectively. Thus, the true parameter vector  $\gamma_0$  of  $\gamma$  can be written as  $\gamma_0 = (\gamma_{10}^\top, \mathbf{0}^\top)^\top$ , where  $\gamma_{10}$  is the true value of  $\gamma_1$ . In addition, the penalized likelihood estimator  $\hat{\gamma}$  of  $\gamma$  can be written as  $\hat{\gamma} = (\hat{\gamma}_1^\top, \hat{\gamma}_2^\top)^\top$ . Let  $\mathcal{A}_\gamma = \{j : \gamma_{0j} \neq 0\}$  be the index set of nonzero components of  $\gamma_0$ , where  $\gamma_{0j}$  is the  $j$ th component of  $\gamma_0$  for  $j = 1, \dots, q$ . Denote the cardinality of  $\mathcal{A}_\gamma$  as  $d_m = |\mathcal{A}_\gamma|$ , which is usually unknown to be estimated in applications. Here, we assume that the non-

---

sparsity size  $d_m \ll n$ , and the dimensionality satisfies  $\log(q) = O(n^\alpha)$  for some  $\alpha \in (0, 1/2)$ . Following Lv and Fan (2009), Zhang (2010), and Fan and Lv (2011), we define the local concavity of the penalty function  $f_{\lambda_n}(t)$  at  $\mathbf{v} = (v_1, \dots, v_q)^\top \in \mathbf{R}^q$  (i.e.,  $\|\mathbf{v}\|_0 = q$ ) as

$$\rho(f_{\lambda_n}; \mathbf{v}) = \lim_{\epsilon \rightarrow 0^+} \max_{1 \leq j \leq q} \sup_{t_1 < t_2 \in (|v_j| - \epsilon, |v_j| + \epsilon)} \frac{\partial_t f_{\lambda_n}(t_2) - \partial_t f_{\lambda_n}(t_1)}{t_2 - t_1},$$

where  $\partial_t^k f_{\lambda_n}(t)$  represents the  $k$ -order derivation of  $f_{\lambda_n}(t)$  with respect to  $t$ , and  $\|\mathbf{A}\|_m$  denotes the  $L_m$  norm of a vector or matrix  $\mathbf{A}$  for  $m \in [0, \infty]$ .

We use  $\text{tr}(\mathbf{A})$  to represent the trace of matrix  $\mathbf{A}$ . Denote  $g_i(\boldsymbol{\gamma}) = \log[\pi(\mathbf{U}_i; \boldsymbol{\gamma}) / \{1 - \pi(\mathbf{U}_i; \boldsymbol{\gamma})\}]$  for  $i = 1, \dots, n$ , and define the Fisher information matrix as

$$\mathbf{F}_n(\boldsymbol{\gamma}) = E\{-\partial^2 l_n(\boldsymbol{\gamma}) / \partial \boldsymbol{\gamma} \partial \boldsymbol{\gamma}^\top\} = \partial_{\boldsymbol{\gamma}}^\top \mathbf{g}(\boldsymbol{\gamma}) \boldsymbol{\Sigma}(\boldsymbol{\gamma}) \partial_{\boldsymbol{\gamma}} \mathbf{g}(\boldsymbol{\gamma}),$$

where  $\mathbf{g}(\boldsymbol{\gamma}) = (g_1(\boldsymbol{\gamma}), \dots, g_n(\boldsymbol{\gamma}))^\top$ , and  $\boldsymbol{\Sigma}(\boldsymbol{\gamma}) = \text{diag}(\pi_1(1 - \pi_1), \dots, \pi_n(1 - \pi_n))$  with  $\pi_i = \pi_i(\boldsymbol{\gamma}) = \pi(\mathbf{U}_i; \boldsymbol{\gamma})$ , for  $i = 1, \dots, n$ . Let  $s_n = \frac{1}{2} \min_j \{|\gamma_{0j}| : \gamma_{0j} \neq 0\}$ , and define  $\mathcal{N} = \{\boldsymbol{\tau} = (\boldsymbol{\tau}_1^\top, \boldsymbol{\tau}_2^\top)^\top \in \mathbf{R}^q : \boldsymbol{\tau}_2 = 0, \|\boldsymbol{\tau}_1 - \boldsymbol{\gamma}_{10}\|_\infty \leq s_n\}$ .

The following assumptions are required to ensure the consistency of  $\hat{\boldsymbol{\gamma}}$ .

**Assumption 4.** The penalty function  $f_{\lambda_n}(t)$  is increasing and concave with respect to  $t \in [0, \infty)$ , and has a continuous derivation  $\partial f_{\lambda_n}(t)$  with  $\partial f_{\lambda_n}(0+) = c_0$ , where  $c_0$  is a positive constant. Also,  $\partial f_{\lambda_n}(t)$  is increasing with respect to  $\lambda_n \in (0, \infty)$ , and  $\partial f_{\lambda_n}(0+)$  is independent of  $\lambda_n$ .

---

Assumption 4 holds for a class of penalty functions, such as the SCAD and MCP penalty functions. By Assumption 4,  $f_{\lambda_n}(t)$  is a concave function with respect to  $t \in [0, \infty)$  to ensure  $\rho(f_{\lambda_n}; \mathbf{v}) \geq 0$ .

**Assumption 5.** (i)  $\min_{\boldsymbol{\tau} \in \mathcal{N}} \mathbb{E}_{\min} \{\mathbf{F}_n(\boldsymbol{\tau})\} \geq cn$  and  $\text{tr}\{\mathbf{F}_n(\boldsymbol{\gamma}_0)\} = O(d_m n)$ ;

(ii)  $\left\| \partial_{\boldsymbol{\gamma}_2}^\top \mathbf{g}(\boldsymbol{\gamma}_0) \boldsymbol{\Sigma}(\boldsymbol{\gamma}_0) \partial_{\boldsymbol{\gamma}_1} \mathbf{g}(\boldsymbol{\gamma}_0) \right\|_{2, \infty} = O(n)$ , where notation  $\|B\|_{2, \infty}$  represents  $\max_{\|\mathbf{v}\|_2=1} \|B\mathbf{v}\|_\infty$ ;

(iii)  $\max_{\boldsymbol{\tau} \in \mathcal{N}, 1 \leq j \leq q} \mathbb{E}_{\max} \left[ \partial_{\boldsymbol{\gamma}_1}^\top \mathbf{g}(\boldsymbol{\tau}) \text{diag}\{|\partial_{\gamma_j} \mathbf{g}(\boldsymbol{\tau})| \circ |\partial_{\boldsymbol{\gamma}_1}^2 \mathbf{V}(\boldsymbol{\tau})|\} \partial_{\boldsymbol{\gamma}_1} \mathbf{g}(\boldsymbol{\tau}) \right] = O(n)$ , where  $\mathbf{V}(\boldsymbol{\gamma}) = (\pi_1(\boldsymbol{\gamma}), \dots, \pi_n(\boldsymbol{\gamma}))^\top$ .

Assumption 5(i) has been used in Fan and Lv (2011), and ensures that the information matrix  $\mathbf{F}_n(\boldsymbol{\tau})$  is positive definite, and its eigenvalues are uniformly bounded. Assumption 5(ii) measures the correlation between each unimportant variable and important variable using the weighted matrix  $\boldsymbol{\Sigma}(\boldsymbol{\gamma}_0)$ , and controls the uniformly growth rate of these regression coefficients. This assumption is similar to the strong irrepresentable condition of Zhao and Yu (2006) for the consistency of Lasso estimator. Assumption 5(iii) is used to control the order of the remainder term when taking the third-order expansion of the objective function.

**Assumption 6.**  $s_n \gg \lambda_n \gg \sqrt{d_m/n}$ ,  $\sqrt{n} \partial f_{\lambda_n}(s_n) = O(1)$  and  $\rho_0 = o(1)$ , where  $\rho_0 = \max_{\boldsymbol{\tau} \in \mathcal{N}} \rho(f_{\lambda_n}; \boldsymbol{\tau})$ .

---

Assumption 6 shows that the minimum signal  $s_n$  should be satisfied, and is used to obtain nice properties of the proposed PLE like other variable selection methods. However, for the  $L_1$  penalty,  $\lambda_n = \partial f_{\lambda_n}(s_n) = O(n^{-1/2})$  is in conflict with the assumption  $\lambda_n \gg \sqrt{d_m/n}$ , which implies that the  $L_1$  penalized likelihood estimator can usually not achieve the consistency rate of  $O_p(\sqrt{d_m/n})$  given in Theorem S1.1, and has not the oracle property like the SCAD penalty function. The assumption  $s_n \gg \lambda_n$  holds automatically for the SCAD penalty function. Thus, Assumption 6 is less restrictive for the SCAD penalty function.

**Theorem S1.1.** *Suppose that Assumptions 4–6 hold. There is a strict local maximizer  $\hat{\gamma} = (\hat{\gamma}_1^\top, \hat{\gamma}_2^\top)^\top$  of the nonconcave penalized likelihood  $Q_n(\gamma)$  with respect to  $\gamma$  such that  $\hat{\gamma}_2 = 0$  with probability tending to 1 as  $n \rightarrow \infty$  and  $\|\hat{\gamma}_1 - \gamma_{10}\|_2 = O_p(\sqrt{d_m/n})$ .*

Theorem S1.1 shows that the sparsity property of the proposed PLE still holds in a high-dimensional parametric model. That is, zero components in  $\gamma_0$  are estimated as zero with probability tending to one. Also, Theorem S1.1 establishes the consistency of the proposed PLE  $\hat{\gamma}_1$  of  $\gamma_1$ , i.e., there is a root- $(n/d_m)$ -consistent PLE of  $\gamma_1$ .

To establish the asymptotic normality of the proposed PLE, we need the following additional assumption, which is associated with the Lyapunov

---

assumptions.

**Assumption 7.**  $\partial f_{\lambda_n}(s_n) = o(1/\sqrt{nd_m})$ ,  $\max_i E|\delta_i - \pi_i(\boldsymbol{\gamma}_0)|^3 = O(1)$ , and  $\sum_{i=1}^n \{\partial_{\boldsymbol{\gamma}_1}^\top g_i(\boldsymbol{\gamma}_0) \mathbf{F}_{n11}^{-1}(\boldsymbol{\gamma}_0) \partial_{\boldsymbol{\gamma}_1} g_i(\boldsymbol{\gamma}_0)\}^{3/2} \rightarrow 0$  as  $n \rightarrow \infty$ , where  $\mathbf{F}_{n11}(\boldsymbol{\gamma}_0) = \partial_{\boldsymbol{\gamma}_1}^\top \mathbf{g}(\boldsymbol{\gamma}_0) \boldsymbol{\Sigma}(\boldsymbol{\gamma}_0) \partial_{\boldsymbol{\gamma}_1} \mathbf{g}(\boldsymbol{\gamma}_0)$ .

**Theorem S1.2.** *Suppose that Assumptions 4–7 hold and  $d_m = o(n^{1/4})$ .*

*Then, we have*

(i) (Sparsity)  $\hat{\boldsymbol{\gamma}}_2 = 0$  with probability tending to 1 as  $n \rightarrow \infty$ .

(ii) (Asymptotic Normality)  $\mathbf{U}_n \mathbf{F}_{n11}^{1/2}(\boldsymbol{\gamma}_0) (\hat{\boldsymbol{\gamma}}_1 - \boldsymbol{\gamma}_{10}) \xrightarrow{\mathcal{L}} \mathcal{N}(\mathbf{0}, \mathbf{G})$ , where  $\mathbf{U}_n$  is an  $m \times d$  matrix such that  $\mathbf{U}_n \mathbf{U}_n^\top \rightarrow \mathbf{G}$ ,  $\mathbf{G}$  is an  $m \times m$  symmetric positive definite matrix with the fixed  $m$ , and  $\xrightarrow{\mathcal{L}}$  represents convergence in distribution.

Theorem S1.2 indicates that the sparsity and asymptotic normality of the proposed PLE still hold even for dimensionality of nonpolynomial order of sample size.

## S2. Properties of the proposed screening procedure

Under the assumption  $Y \perp \delta | X_k$  for  $k = 1, \dots, p$ ,  $\hat{r}_k$  can be regarded as the empirical estimator of  $E(X_k Y)$  in the presence of responses MAR. Without loss of generality, for each  $k = 1, \dots, p$ , the  $k$ th column of covariates satisfies  $E(X_k) = 0$ ,  $E(X_k^2) = 1$ . Then we have  $E(X_k Y) = \text{cov}(X_k, Y)$

---

and  $\beta_k = \text{cov}(X_k, Y)$ , which indicates that  $\beta_k$  is the covariance between  $X_k$  and  $Y$ . Hence,  $\beta_k = 0$  is equivalent to the fact that  $X_k$  and  $Y$  are marginally uncorrelated. Thus, define the index set of the active predictors as  $\mathcal{M}_* = \{k : \beta_k \neq 0 \text{ for } 1 \leq k \leq p\}$ , which corresponds to the true sparse model with nonsparsity size  $|\mathcal{M}_*|$ , where  $|\mathcal{M}_*|$  is the cardinality of  $\mathcal{M}_*$ , and denote  $\mathcal{I}_* = \{1, \dots, p\} \setminus \mathcal{M}_*$  as the index set of the inactive predictors, where  $p \gg n$ . Here, we assume that  $p \gg |\mathcal{M}_*|$  in ultrahigh-dimensional data analysis and define  $r_k = E(X_k Y)$ . To investigate the sure screening properties of the presented screening criterion, we require the following assumptions.

**Assumption 8.** For  $k = 1, \dots, p$ , the probability density function of  $X_k$ , say  $f_k(x)$ , has continuous and bounded second order derivatives over the support  $\mathcal{X}_k$  of  $X_k$ , and is bounded away from zero and infinity uniformly over  $\mathcal{X}_k$ .

**Assumption 9.** The kernel function  $K(\cdot)$  is a probability density function such that (i) it is bounded and has compact support; (ii) it is symmetric with  $\int t^2 K(t) dt < 1$ ; (iii)  $K(\cdot) \geq d_1$  for some positive constant  $d_1$  in some closed interval centered at zero; (iv)  $\sqrt{nh^2} \rightarrow 0$  as  $n \rightarrow \infty$ .

**Assumption 10.** Variables  $X_k$ ,  $Y$  and  $X_k Y$  satisfy the sub-exponential tail probability uniformly in  $p$ . That is, there exists a positive constant  $u_0$

---

such that for all  $0 < \tilde{u} \leq 2u_0$ ,

$$\max_{1 \leq k \leq p} E\{\exp(2\tilde{u}X_k^2)\} < \infty, \quad E\{\exp(2\tilde{u}Y^2)\} < \infty,$$

$$\max_{1 \leq k \leq p} E\{\exp(2\tilde{u}X_kY)\} < \infty.$$

**Assumption 11.** There exists a positive constant  $c_0 > 0$  and  $0 < \varsigma < 1/2$  such that  $\min_{k \in \mathcal{M}_*} |r_k| \geq c_0 n^{-\varsigma}$ .

**Assumption 12.**  $\lim_{p \rightarrow \infty} \inf\{\min_{k \in \mathcal{M}_*} |r_k| - \max_{k \notin \mathcal{M}_*} |r_k|\} \geq m_0$  for some  $m_0 > 0$ .

Assumptions 8 and 9 impose some regularity assumptions on the probability density functions  $f_k(x)$  and kernel function  $K(\cdot)$ , respectively, which hold for the widely used distributions and kernel functions. The assumption that  $\sqrt{nh^2} \rightarrow 0$  is to control the bias induced by the kernel smoothing. Assumption 10 has widely used in high-dimensional data analysis (Fan and Lv, 2008, and Li et al., 2012), and holds if  $\mathbf{X}$ ,  $Y$  and  $X_kY$  are bounded uniformly of  $\mathbf{X}$ ,  $Y$ , and  $X_kY$  have multivariate normal distribution. Assumption 11 allows the minimal signal between active variables and response variable to be the order of  $n^{-\varsigma}$ , which is a widely used condition to guarantee the sure screening property. Assumption 12 ensures that the active and inactive predictors can be well separated in the population level. Assumption 12 is similar to Condition (C3) of Cui et al. (2015).

---

**Theorem S2.1.** (*Sure Screening Property*) Under Assumptions 1 and 8–10, then for any constant  $c_2 > 0$ , there exists  $c_7 > 0$  such that

$$\Pr\left(\max_{1 \leq k \leq p} |\hat{r}_k - r_k| \geq c_2 n^{-\varsigma}\right) \leq O\{p \exp(-c_7 n^{(1-2\varsigma)/3} + \log(n))\}$$

for  $n$  sufficiently large. Furthermore, under Assumption 11, by taking  $\varrho_n = c_8 n^{-\varsigma}$  with  $c_8 \leq c_0/2$ , there exists some positive constant  $c_9$  such that

$$\Pr(\mathcal{M}_* \subset \widehat{\mathcal{M}}) \geq 1 - O\{|\mathcal{M}_*| \exp(-c_9 n^{(1-2\varsigma)/3} + \log(n))\}.$$

(*Ranking Consistency Property*) If Assumptions 1, 8–12 and additional assumptions  $\log(p) = o(n^{1/3} m_0^{2/3})$  and  $\log(n) = o(n^{1/3} m_0^{2/3})$  hold. Then, we have

$$\liminf_{n \rightarrow \infty} \left\{ \min_{k \in \mathcal{M}_*} |\hat{r}_k| - \max_{k \in \mathcal{I}_*} |\hat{r}_k| \right\} > 0, a.s..$$

Theorem S2.1 shows the sure screening and rank consistency properties of the proposed screening procedure, which indicates that the proposed MI-SIS method can handle the NP-dimensionality problem. Specifically, as  $n \rightarrow \infty$ , the maximum dimensional is  $p = o\{\exp(n^{(1-2\varsigma)/3})\}$  for  $\varsigma \in (0, 1/2)$ . In addition, through the ranking consistency property, we can separate the active and inactive predictors by taking an ideal threshold value  $\varrho_n$ , which is smaller than the minimum signal.



---

### S3. Proofs of all theorems

Let  $\mathbf{P}_s^* = \mathbf{D}_s(\mathbf{P}_s - \mathbf{I}) + \mathbf{I}$ , where  $\mathbf{P}_s = \mathbf{X}_s(\mathbf{X}_s^\top \mathbf{X}_s)^{-1} \mathbf{X}_s$ ,  $\mathbf{D}_s = \text{diag}(d_1^s, \dots, d_n^s)$ ,

$d_i^s = 1/(1-h_{ii}^s)$  and  $h_{ii}^s$  is the  $i$ th diagonal element of  $\mathbf{P}_s = \mathbf{X}_s(\mathbf{X}_s^\top \mathbf{X}_s)^{-1} \mathbf{X}_s^\top$ .

Let  $\|\mathbf{A}\|$  be the Frobenius norm (e.g.,  $\|\mathbf{A}\| = \sqrt{\text{tr}(\mathbf{A}^\top \mathbf{A})}$ ), and  $\|\mathbf{A}\|_2$  be the

Euclidean norm (e.g.,  $\|\mathbf{A}\|_2 = \sqrt{\mathbb{E}_{\max}(\mathbf{A}^* \mathbf{A})}$ ) of matrix  $\mathbf{A}$ , where  $\mathbf{A}^*$  rep-

resents the conjugate transpose of matrix  $\mathbf{A}$ . Denote  $\mathbf{P}^*(\boldsymbol{\omega}) = \sum_{s=1}^S \omega_s \mathbf{P}_s^*$

and  $\mathbf{P}(\boldsymbol{\omega}) = \sum_{s=1}^S \omega_s \mathbf{P}_s$ .

**Proof of Theorem S1.1.** First, we show the consistency of the proposed

PLE in the  $d_m$ -dimensional subspace. To this end, we constrain the  $Q_n(\boldsymbol{\gamma})$

on the  $d_m$ -dimensional subspace  $\{\boldsymbol{\gamma} \in \mathbb{R}^q : \gamma_j = 0, j \in \mathcal{A}_\gamma^c\}$  of  $\mathbb{R}^q$ , and the

corresponding constrained penalized likelihood function is given by

$$\tilde{Q}_n(\boldsymbol{\tau}) = \frac{1}{n} \tilde{l}_n(\boldsymbol{\tau}) - \sum_{j=1}^{d_m} f_{\lambda_n}(|\tau_j|), \quad (\text{S3.1})$$

where  $\boldsymbol{\tau} = (\tau_1, \dots, \tau_{d_m})^\top$  and  $\tilde{l}_n(\boldsymbol{\tau}) = \sum_{i=1}^n [\delta_i \log\{\pi_i(\mathbf{U}_i^\tau; \boldsymbol{\tau})\} + (1-\delta_i) \log\{1 - \pi(\mathbf{U}_i^\tau; \boldsymbol{\tau})\}]$  in which  $\mathbf{U}_i^\tau$  is the subvector of  $\mathbf{U}_i$  corresponding to  $\boldsymbol{\tau}$ .

Let  $\alpha_n = \sqrt{d_m/n}$ , and define the closed set  $\mathcal{N}_0 = \{\boldsymbol{\tau} \in \mathbb{R}^{d_m} : \|\boldsymbol{\tau} - \boldsymbol{\gamma}_{10}\|_2 \leq \alpha_n u\}$  for  $u \in (0, \infty)$ . Here, the purpose is to show that, for any

$\kappa > 0$  and a sufficiently large  $n$ , we have

$$\Pr \left\{ \sup_{\boldsymbol{\tau} \in \partial \mathcal{N}_0} \tilde{Q}_n(\boldsymbol{\tau}) < \tilde{Q}_n(\boldsymbol{\gamma}_{10}) \right\} \geq 1 - \kappa,$$

where  $\partial \mathcal{N}_0$  denotes the boundary of the closed set  $\mathcal{N}_0$ , which implies that

---

there exists a local maximizer  $\hat{\gamma}_1$  in  $\mathcal{N}_0$  such that  $\|\hat{\gamma}_1 - \gamma_{10}\|_2 = O_p(\sqrt{d_m/n})$ .

For a sufficiently large  $n$ , it follows from Assumption 6 that  $\alpha_n u \leq d_m$ .

Taking Taylor expansion of  $\tilde{Q}_n(\boldsymbol{\tau})$  at the true value  $\gamma_{10}$  of  $\gamma_1$  yields

$$\tilde{Q}_n(\boldsymbol{\tau}) - \tilde{Q}_n(\gamma_{10}) = (\boldsymbol{\tau} - \gamma_{10})^\top \mathbf{D}_1 - \frac{1}{2}(\boldsymbol{\tau} - \gamma_{10})^\top \mathbf{D}_2(\boldsymbol{\tau} - \gamma_{10}) \quad (\text{S3.2})$$

for any  $\boldsymbol{\tau} \in \mathcal{N}_0$ , where  $\mathbf{D}_1 = \partial \tilde{Q}_n(\gamma_{10}) / \partial \gamma_1 = \partial_{\gamma_1}^\top \mathbf{g}(\gamma_{10}) \{\boldsymbol{\delta} - \mathbf{V}(\gamma_{10})\} / n - \partial f_{\lambda_n}(\gamma_{10})$ ,  $\mathbf{D}_2 = \partial^2 \tilde{Q}_n(\tilde{\gamma}_1) / \partial \gamma_1 \partial \gamma_1^\top$ , and  $\tilde{\gamma}_1$  lies on the line segment jointing  $\boldsymbol{\tau}$  and  $\gamma_{10}$ . Following Fan and Peng (2004), we can obtain

$$\begin{aligned} \mathbf{D}_2 &= -\frac{1}{n} \partial^2 \tilde{l}_n(\tilde{\gamma}_1) / \partial \gamma_1 \partial \gamma_1^\top + \text{diag}\{\partial^2 f_{\lambda_n}(|\tilde{\gamma}_1|)\} \\ &= \frac{1}{n} \mathbf{F}_n(\tilde{\gamma}_1) - \frac{1}{n} \left\{ \partial^2 \tilde{l}_n(\tilde{\gamma}_1) / \partial \gamma_1 \partial \gamma_1^\top - E(\partial^2 \tilde{l}_n(\tilde{\gamma}_1) / \partial \gamma_1 \partial \gamma_1^\top) \right\} \\ &\quad + \text{diag}\{\partial^2 f_{\lambda_n}(|\tilde{\gamma}_1|)\} \\ &= \frac{1}{n} \mathbf{F}_n(\tilde{\gamma}_1) + \text{diag}\{\partial^2 f_{\lambda_n}(|\tilde{\gamma}_1|)\} + o_p(1). \end{aligned}$$

Without loss of generality, when there is no the second-order derivative of the penalty function  $f_{\lambda_n}(\cdot)$ , it is easily shown that matrix  $\mathbf{D}_2$  can be replaced by a diagonal matrix whose maximum absolute element is bounded by  $\rho_0$ . Thus, for a sufficiently large  $n$ , under Assumptions 5(i) and 6, we have  $\mathbb{E}_{\min}(\mathbf{D}_2) \geq c - \rho_0 \geq \frac{c}{2}$ , where  $c$  is a constant.

Using the Markov's inequality and (S3.2) yields

$$\Pr \left\{ \sup_{\boldsymbol{\tau} \in \mathcal{N}_0} \tilde{Q}_n(\boldsymbol{\tau}) < \tilde{Q}_n(\gamma_{10}) \right\} \geq \Pr \left\{ u \alpha_n \left( \|\mathbf{D}_1\|_2 - \frac{cu \alpha_n}{4} \right) < 0 \right\} \geq 1 - \frac{16E\|\mathbf{D}_1\|_2^2}{c^2 u^2 \alpha_n^2}.$$

---

However, under Assumptions 5(i) and 6, we have

$$E\|\mathbf{D}_1\|_2^2 \leq \frac{1}{n^2} \text{tr}\{\mathbf{F}_n(\boldsymbol{\gamma}_{10})\} + d_m \partial f_{\lambda_n}(d_m)^2 = O(\alpha_n^2),$$

which leads to

$$\Pr \left\{ \sup_{\boldsymbol{\tau} \in \partial \mathcal{N}_0} \tilde{Q}_n(\boldsymbol{\tau}) < \tilde{Q}_n(\boldsymbol{\gamma}_{10}) \right\} \geq 1 - O\left(\frac{16}{c^2 u^2}\right).$$

Thus, we prove  $\|\hat{\boldsymbol{\gamma}}_1 - \boldsymbol{\gamma}_{10}\|_2 = O_p(\sqrt{d_m/n})$ .

Next, to show the sparsity of the proposed estimator, it is necessary to prove that  $\hat{\boldsymbol{\gamma}} \in R^q$  is a strict local maximizer of  $Q_n(\boldsymbol{\gamma})$  such that  $\hat{\boldsymbol{\gamma}}_{\mathcal{A}_\gamma} = \hat{\boldsymbol{\gamma}}_1 \in \mathcal{N}_0 \subset \mathcal{N}$  and  $\hat{\boldsymbol{\gamma}}_{\mathcal{A}_\gamma^c} = \hat{\boldsymbol{\gamma}}_2 = 0$ . Similar to the proof of Theorem 1 of Fan and Lv (2011), we only require showing

$$\left\| \frac{1}{n\lambda_n} \partial_{\boldsymbol{\gamma}_2} l_n(\hat{\boldsymbol{\gamma}}) \right\|_\infty \leq \partial f_{\lambda_n}(0+). \quad (\text{S3.3})$$

Thus, we have

$$\begin{aligned} \frac{1}{n} \partial_{\boldsymbol{\gamma}_2} l_n(\hat{\boldsymbol{\gamma}}) &= \frac{1}{n} \partial_{\boldsymbol{\gamma}_2}^\top \mathbf{g}(\hat{\boldsymbol{\gamma}}) \{\boldsymbol{\delta} - \mathbf{V}(\hat{\boldsymbol{\gamma}})\} \\ &= \frac{1}{n} [\boldsymbol{\eta}_{\mathcal{A}_\gamma^c} + \{\partial_{\boldsymbol{\gamma}_2}^\top \mathbf{g}(\hat{\boldsymbol{\gamma}}) \boldsymbol{\delta} - \partial_{\boldsymbol{\gamma}_2}^\top \mathbf{g}(\hat{\boldsymbol{\gamma}}) \mathbf{V}(\hat{\boldsymbol{\gamma}})\} \\ &\quad - \{\partial_{\boldsymbol{\gamma}_2}^\top \mathbf{g}(\boldsymbol{\gamma}_0) \boldsymbol{\delta} - \partial_{\boldsymbol{\gamma}_2}^\top \mathbf{g}(\boldsymbol{\gamma}_0) \mathbf{V}(\boldsymbol{\gamma}_0)\}], \end{aligned}$$

where  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_q)^\top = \partial_{\boldsymbol{\gamma}}^\top \mathbf{g}(\boldsymbol{\gamma}_0) \{\boldsymbol{\delta} - \mathbf{V}(\boldsymbol{\gamma}_0)\}$ .

Denote  $\mathcal{B}_\gamma = \{\|\boldsymbol{\eta}_{\mathcal{A}_\gamma^c}\|_\infty \leq c_1 \sqrt{n}\}$ , where  $c_1$  is some constant. Under Assumption 1, it is easily known that the first- to third-order derivatives of  $g_i(\boldsymbol{\gamma})$  with respect to  $\boldsymbol{\gamma}$  are bounded. By Bonferroni's inequality and

---

$\log(q) = O(n^\alpha)$  for some  $\alpha \in (0, 1/2)$ , there exists a constant  $c'$  such that

$$\begin{aligned} \Pr(\mathcal{B}_\gamma) &= \Pr(\|\boldsymbol{\eta}_{\mathcal{A}_\gamma^c}\|_\infty \leq c_1\sqrt{n}) = 1 - \Pr(\|\boldsymbol{\eta}_{\mathcal{A}_\gamma^c}\|_\infty > c_1\sqrt{n}) \\ &\geq 1 - 2(p-d)\exp(-c'n) \geq 1 - 2q\exp(-c'n) \rightarrow 1. \end{aligned}$$

Under Assumptions 5(ii), 5(iii) and 6, we consider the Taylor expansion of  $\partial_{\boldsymbol{\gamma}_2}^\top \mathbf{g}(\hat{\boldsymbol{\gamma}})\boldsymbol{\delta} - \partial_{\boldsymbol{\gamma}_2}^\top \mathbf{g}(\hat{\boldsymbol{\gamma}})\mathbf{V}(\hat{\boldsymbol{\gamma}})$  at  $\boldsymbol{\gamma}_{10}$ . Thus, we obtain

$$\begin{aligned} \left\| \frac{1}{n\lambda_n} \partial_{\boldsymbol{\gamma}_2} l_n(\hat{\boldsymbol{\gamma}}) \right\|_\infty &\leq \frac{1}{n\lambda_n} \left[ \|\boldsymbol{\eta}_{\mathcal{A}_\gamma^c}\|_\infty + \left\| \left\{ \partial_{\boldsymbol{\gamma}_2}^\top \mathbf{g}(\hat{\boldsymbol{\gamma}})\boldsymbol{\delta} - \partial_{\boldsymbol{\gamma}_2}^\top \mathbf{g}(\hat{\boldsymbol{\gamma}})\mathbf{V}(\hat{\boldsymbol{\gamma}}) \right\} \right. \\ &\quad \left. - \left\{ \partial_{\boldsymbol{\gamma}_2}^\top \mathbf{g}(\boldsymbol{\gamma}_0)\boldsymbol{\delta} - \partial_{\boldsymbol{\gamma}_2}^\top \mathbf{g}(\boldsymbol{\gamma}_0)\mathbf{V}(\boldsymbol{\gamma}_0) \right\} \right\|_\infty \\ &\leq o_p(1) + \frac{1}{n\lambda_n} \left\{ O(n)\|\hat{\boldsymbol{\gamma}}_1 - \boldsymbol{\gamma}_{10}\|_2 + O(n)\|\hat{\boldsymbol{\gamma}}_1 - \boldsymbol{\gamma}_{10}\|_2^2 \right\} \\ &= o_p(1) + O_p(\lambda_n^{-1}\sqrt{d_m/n}) = o_p(1). \end{aligned}$$

Then, for a sufficiently large  $n$ , (S3.3) holds. Hence, we finish the proof of Theorem S1.1.

**Proof of Theorem S1.2.** From the proof of Theorem S1.1, we only require proving the asymptotic normality of  $\hat{\boldsymbol{\gamma}}_1$ . For the set  $\mathcal{N}$ , it is easily shown that  $\hat{\boldsymbol{\gamma}} = (\hat{\boldsymbol{\gamma}}_1^\top, \mathbf{0}^\top)^\top \in \mathcal{N}$  is a strict local maximizer of  $\tilde{Q}_n(\boldsymbol{\tau})$ , which leads to  $\partial_{\boldsymbol{\gamma}_1} \tilde{Q}_n(\hat{\boldsymbol{\gamma}}) = 0$ .

Taking Taylor expansion of  $\partial_{\boldsymbol{\gamma}_1} \tilde{l}_n(\hat{\boldsymbol{\gamma}})$  at  $\boldsymbol{\gamma}_{10}$  leads to

$$\begin{aligned} 0 &= \partial_{\boldsymbol{\gamma}_1} \tilde{Q}_n(\hat{\boldsymbol{\gamma}}) = \frac{1}{n} \partial_{\boldsymbol{\gamma}_1} \tilde{l}_n(\hat{\boldsymbol{\gamma}}) - \partial f_{\lambda_n}(\hat{\boldsymbol{\gamma}}) \\ &= \frac{1}{n} \partial_{\boldsymbol{\gamma}_1}^\top \mathbf{g}(\boldsymbol{\gamma}_0) \{ \boldsymbol{\delta} - \mathbf{V}(\boldsymbol{\gamma}_0) \} \\ &\quad - \frac{1}{n} \left[ \mathbf{F}_{n11}(\boldsymbol{\gamma}_0) - \left\{ \frac{\partial^2 l_n(\boldsymbol{\gamma}_0)}{\partial \boldsymbol{\gamma}_1 \partial \boldsymbol{\gamma}_1^\top} - E \left( \frac{\partial^2 l_n(\boldsymbol{\gamma}_0)}{\partial \boldsymbol{\gamma}_1 \partial \boldsymbol{\gamma}_1^\top} \right) \right\} \right] (\hat{\boldsymbol{\gamma}}_1 - \boldsymbol{\gamma}_{10}) \end{aligned}$$

---


$$\begin{aligned}
& + O(1)\|\hat{\gamma}_1 - \gamma_{10}\|_2^2 d_m - \partial f_{\lambda_n}(\hat{\gamma}) \\
& = \frac{1}{n} \partial_{\gamma_1}^\top \mathbf{g}(\gamma_0) \{\boldsymbol{\delta} - \mathbf{V}(\gamma_0)\} - \frac{1}{n} \mathbf{F}_{n11}(\gamma_0) (\hat{\gamma}_1 - \gamma_{10}) - \partial f_{\lambda_n}(\hat{\gamma}) \\
& + o_p(n^{-1/2} d_m^{-1/2}) + O_p(d_m^2 n^{-1}).
\end{aligned}$$

Under Assumption 7, it follows from  $\min_j |\tau_j| \geq \min_{j \in \mathcal{A}_\gamma} |\gamma_{0j}| - s_n = s_n$  and the monotonicity of  $\partial f_{\lambda_n}(t)$  that

$$\|\partial f_{\lambda_n}(\hat{\gamma})\|_2 \leq \sqrt{s_n} \partial f_{\lambda_n}(s_n) = o_p(n^{-1/2}). \quad (\text{S3.4})$$

Combining  $d_m = o(n^{1/4})$ , (S2.4) and (S3.4) leads to

$$\mathbf{F}_{n11}(\gamma_0) (\hat{\gamma}_1 - \gamma_{10}) = \partial_{\gamma_1}^\top \mathbf{g}(\gamma_0) \{\boldsymbol{\delta} - \mathbf{V}(\gamma_0)\} + o_p(\sqrt{n}).$$

By Assumption 5(i), we obtain

$$\mathbf{F}_{n11}^{1/2}(\gamma_0) (\hat{\gamma}_1 - \gamma_{10}) = \mathbf{F}_{n11}^{-1/2}(\gamma_0) \partial_{\gamma_1}^\top \mathbf{g}(\gamma_0) \{\boldsymbol{\delta} - \mathbf{V}(\gamma_0)\} + o_p(1). \quad (\text{S3.5})$$

Let  $\mathbf{U}_n \mathbf{U}_n^\top = \mathbf{G}$ , where  $\mathbf{U}_n$  is an  $m \times d_m$  matrix, and  $\mathbf{G}$  is an  $m \times m$  symmetric positive definite matrix. Then, we have  $\mathbf{U}_n \mathbf{F}_{n11}^{1/2}(\gamma_0) (\hat{\gamma}_1 - \gamma_{10}) = \boldsymbol{\nu}_n + o_p(1)$ , where  $\boldsymbol{\nu}_n = \mathbf{U}_n \mathbf{F}_{n11}^{-1/2}(\gamma_0) \partial_{\gamma_1}^\top \mathbf{g}(\gamma_0) \{\boldsymbol{\delta} - \mathbf{V}(\gamma_0)\}$ .

If we can prove  $\boldsymbol{\nu}_n \xrightarrow{\mathcal{L}} \mathcal{N}(\mathbf{0}, \mathbf{G})$ , it follows from Slutsky's theorem that  $\mathbf{U}_n \mathbf{F}_{n11}^{1/2}(\gamma_0) (\hat{\gamma}_1 - \gamma_{10}) \xrightarrow{\mathcal{L}} \mathcal{N}(\mathbf{0}, \mathbf{G})$ . For any unit vector  $\mathbf{a} \in \mathbf{R}^m$ , we have

$$\mathbf{v}_n = \mathbf{a}^\top \boldsymbol{\nu}_n = \mathbf{a}^\top \mathbf{U}_n \mathbf{F}_{n11}^{-1/2}(\gamma_0) \partial_{\gamma_1}^\top \mathbf{g}(\gamma_0) \{\boldsymbol{\delta} - \mathbf{V}(\gamma_0)\} = \sum_{i=1}^n z_i,$$

where  $z_i = \mathbf{a}^\top \mathbf{U}_n \mathbf{F}_{n11}^{-1/2}(\gamma_0) \partial_{\gamma_1}^\top \mathbf{g}_i(\gamma_0) \{\delta_i - V_i(\gamma_0)\}$ . It is easily seen that  $z_i$ 's

are independent. For  $i = 1, \dots, n$ , we have  $E(z_i) = 0$  and

$$\sum_{i=1}^n \text{var}(z_i) = \mathbf{a}^\top \mathbf{U}_n \mathbf{F}_{n11}^{-1/2}(\gamma_0) \mathbf{F}_{n11}(\gamma_0) \mathbf{F}_{n11}^{-1/2}(\gamma_0) \mathbf{U}_n^\top \mathbf{a} = \mathbf{a}^\top \mathbf{U}_n \mathbf{U}_n^\top \mathbf{a} \xrightarrow{\mathcal{P}} \mathbf{a}^\top \mathbf{G} \mathbf{a},$$

where  $\xrightarrow{\mathcal{P}}$  represents the convergence in probability.

Under Assumption 7, it follows from Cauchy-Schwarz inequality that

$$\begin{aligned} \sum_{i=1}^n E|z_i|^3 &= \sum_{i=1}^n |\mathbf{a}^\top \mathbf{U}_n \mathbf{F}_{n11}^{-1/2}(\gamma_0) \partial_{\gamma_1} g_i(\gamma_0)|^3 E|\delta_i - V_i(\gamma_0)|^3 \\ &= O(1) \sum_{i=1}^n |\mathbf{a}^\top \mathbf{U}_n \mathbf{F}_{n11}^{-1/2}(\gamma_0) \partial_{\gamma_1} g_i(\gamma_0)|^3 \\ &\leq O(1) \sum_{i=1}^n \|\mathbf{a}^\top \mathbf{U}_n\|_2^3 \cdot \|\mathbf{F}_{n11}^{-1/2}(\gamma_0) \partial_{\gamma_1} g_i(\gamma_0)\|_2^3 \\ &= O(1) \sum_{i=1}^n \{\partial_{\gamma_1} g_i(\gamma_0) \mathbf{F}_{n11}^{-1}(\gamma_0) \partial_{\gamma_1} g_i(\gamma_0)\}^{3/2} = o(1). \end{aligned}$$

Using the Lyapunov's Theorem leads to  $\mathbf{a}^\top \boldsymbol{\nu}_n = \sum_{i=1}^n z_i \xrightarrow{\mathcal{L}} \mathcal{N}(\mathbf{0}, \mathbf{a}^\top \mathbf{G} \mathbf{a})$ , which holds for any unit vector  $\mathbf{a} \in \mathbb{R}^m$ . Thus, we finish the proof of Theorem S1.2(ii).

**Lemma 1.** *Suppose that Assumptions 1, 2 and 3(ii) hold. Let  $p_m = \sup_{1 \leq s \leq S} p_s$ . Then, as  $n \rightarrow \infty$ , we have*

- (i)  $\sup_{1 \leq s \leq S} \left\| \frac{1}{n} \mathbf{X}_s^\top \widehat{\mathbf{W}} \mathbf{X}_s - \frac{1}{n} \mathbf{X}_s^\top \mathbf{W} \mathbf{X}_s \right\| = O_p \left( p_m \sqrt{\frac{d_m}{n}} \right);$
- (ii)  $\sup_{1 \leq s \leq S} \left\| \left( \frac{1}{n} \mathbf{X}_s^\top \widehat{\mathbf{W}} \mathbf{X}_s \right)^{-1} - \left( \frac{1}{n} \mathbf{X}_s^\top \mathbf{W} \mathbf{X}_s \right)^{-1} \right\|_2 = O_p \left( p_m \sqrt{\frac{d_m}{n}} \right);$
- (iii)  $\sup_{1 \leq s \leq S} \left\| \left( \frac{1}{n} \mathbf{X}_s^\top \widehat{\mathbf{W}} \mathbf{X}_s \right)^{-1} \right\|_2 = O_p(1);$
- (iv)  $\widetilde{\mathbf{P}}(\boldsymbol{\omega}) - \widehat{\mathbf{P}}(\boldsymbol{\omega}) = \mathbf{P}^*(\boldsymbol{\omega}) - \mathbf{P}(\boldsymbol{\omega}) + O_p \left( p_m^2 S \sqrt{\frac{d_m}{n}} \right).$

---

**Proof.** By the definition of  $\widehat{\mathbf{W}}$  and  $\mathbf{W}$ , we obtain

$$\begin{aligned} \left\| \frac{1}{n} \mathbf{X}_s^\top \widehat{\mathbf{W}} \mathbf{X}_s - \frac{1}{n} \mathbf{X}_s^\top \mathbf{W} \mathbf{X}_s \right\| &= \left\| n^{-1} \sum_{i=1}^n \left( \frac{\delta_i}{\widehat{\pi}_i} - \frac{\delta_i}{\pi_i} \right) X_{si} X_{si}^\top \right\| \\ &\leq \sup_i \left| \frac{\delta_i}{\widehat{\pi}_i} - \frac{\delta_i}{\pi_i} \right| n^{-1} \sum_{i=1}^n \|X_{si} X_{si}^\top\|. \end{aligned}$$

Following Lemma 1 of Hirano et al. (2003), and by Theorem S1.1 and Assumption 1, we have

$$\sup_i \left| \frac{1}{\widehat{\pi}_i} - \frac{1}{\pi_i} \right| \leq \sup_i \left| \frac{\partial \pi_i / \partial \gamma_1}{\pi_i^2} \right| \cdot \|\widehat{\gamma}_1 - \gamma_{01}\|_2 = O_p \left( \sqrt{\frac{d_m}{n}} \right).$$

From Assumption 2, we have  $\sup_s \sum_{i=1}^n \|X_{si} X_{si}^\top\| / n = O_p(p_m)$ . Thus, we have proved (i). For (ii), let  $\widehat{\mathbf{H}}_s = \mathbf{X}_s^\top \widehat{\mathbf{W}} \mathbf{X}_s / n$  and  $\mathbf{H}_s = \mathbf{X}_s^\top \mathbf{W} \mathbf{X}_s / n$ .

Following Theorem 1 of Lewis and Reinsel (1985), we obtain

$$\widehat{\mathbf{H}}_s^{-1} - \mathbf{H}_s^{-1} = -\widehat{\mathbf{H}}_s^{-1} (\widehat{\mathbf{H}}_s - \mathbf{H}_s) \mathbf{H}_s^{-1} = -\{\mathbf{H}_s^{-1} + (\widehat{\mathbf{H}}_s^{-1} - \mathbf{H}_s^{-1})\} (\widehat{\mathbf{H}}_s - \mathbf{H}_s) \mathbf{H}_s^{-1}.$$

Assumptions 1 and 2 imply that  $\sup_s \|\mathbf{H}_s^{-1}\|_2 \leq \Lambda < \infty$  for some constant  $\Lambda$ . According to Assumption 3(ii) and Lemma 1(i), and  $\sup_s \|\widehat{\mathbf{H}}_s - \mathbf{H}_s\|_2 \leq \sup_s \|\widehat{\mathbf{H}}_s - \mathbf{H}_s\| \rightarrow 0$ , there exists a constant  $\Lambda'$  such that  $\Lambda \|\widehat{\mathbf{H}}_s - \mathbf{H}_s\|_2 < \Lambda' < 1$  as  $n \rightarrow \infty$ . Then, with the probability tending to one, we have

$$\sup_{1 \leq s \leq S} \|\widehat{\mathbf{H}}_s^{-1} - \mathbf{H}_s^{-1}\|_2 \leq \sup_{1 \leq s \leq S} \frac{\Lambda^2 \|\widehat{\mathbf{H}}_s - \mathbf{H}_s\|_2}{1 - \Lambda \|\widehat{\mathbf{H}}_s - \mathbf{H}_s\|_2},$$

which leads to (ii). Using the Triangle inequality and Assumption 3(ii) yields

$$\sup_{1 \leq s \leq S} \|\widehat{\mathbf{H}}_s^{-1}\|_2 \leq \sup_s \|\mathbf{H}_s^{-1}\|_2 + \sup_{1 \leq s \leq S} \|\widehat{\mathbf{H}}_s^{-1} - \mathbf{H}_s^{-1}\|_2 = O_p(1),$$

---

which shows that Lemma 1(iii) holds. It follows from Wiener and Masani (1958) that  $\|\mathbf{A}\mathbf{B}\| \leq \|\mathbf{A}\|_2\|\mathbf{B}\| + \|\mathbf{A}\|\|\mathbf{B}\|_2$  for any matrices  $\mathbf{A}$  and  $\mathbf{B}$ .

Combining Theorem S1.1 and Lemma 1(iii) leads to

$$\begin{aligned}
\sup_{1 \leq s \leq S} \|\widehat{\mathbf{P}}_s - \mathbf{P}_s\| &= \sup_{1 \leq s \leq S} \left\| \frac{1}{n} \mathbf{X}_s \widehat{\mathbf{H}}_s^{-1} \mathbf{X}_s^\top \widehat{\mathbf{W}} - \frac{1}{n} \mathbf{X}_s (\frac{1}{n} \mathbf{X}_s^\top \mathbf{X}_s)^{-1} \mathbf{X}_s^\top \right\| \\
&\leq \sup_i \left| \frac{\delta_i}{\hat{\pi}_i} - \frac{\delta_i}{\pi_i} \right| \cdot \sup_s \|\widehat{\mathbf{H}}_s^{-1}\|_2 \cdot \sup_s \frac{1}{n} \sum_{i=1}^n \|X_{si} X_{si}^\top\| \\
&\quad + \sup_i \left| \frac{\delta_i}{\pi_i} - 1 \right| \cdot \sup_s \|\widehat{\mathbf{H}}_s^{-1}\|_2 \cdot \sup_s \frac{1}{n} \sum_{i=1}^n \|X_{si} X_{si}^\top\| \\
&\quad + \sup_s \left\| \widehat{\mathbf{H}}_s^{-1} - (\frac{1}{n} \mathbf{X}_s^\top \mathbf{X}_s)^{-1} \right\|_2 \cdot \sup_s \frac{1}{n} \sum_{i=1}^n \|X_{si} X_{si}^\top\| \\
&= O_p \left( p_m^2 \sqrt{\frac{d_m}{n}} \right).
\end{aligned}$$

Similarly, we can obtain  $\sup_s \|\widehat{\mathbf{D}}_s - \mathbf{D}_s\| = O_p \left( p_m^2 \sqrt{\frac{d_m}{n}} \right)$ . It follows from  $\widetilde{\mathbf{P}}_s = \widehat{\mathbf{D}}_s(\widehat{\mathbf{P}}_s - \mathbf{I}) + \mathbf{I}$  and the definition of  $\mathbf{P}_s^*$  that

$$\begin{aligned}
\sup_{1 \leq s \leq S} \|\widetilde{\mathbf{P}}_s - \mathbf{P}_s^*\| &= \sup_{1 \leq s \leq S} \left\| \widehat{\mathbf{D}}_s(\widehat{\mathbf{P}}_s - \mathbf{P}_s) + (\widehat{\mathbf{D}}_s - \mathbf{D}_s)\mathbf{P}_s + (\mathbf{D}_s - \widehat{\mathbf{D}}_s) \right\| \\
&= O_p \left( p_m^2 \sqrt{\frac{d_m}{n}} \right),
\end{aligned}$$

which shows that (iv) holds.

**Lemma 2.** *Under Assumption 3(v), there exists a constant  $C > 0$  such that (i)  $\mathbb{E}_{\max}(\mathbf{P}_s^* - \mathbf{P}_s) \leq Cp_s/n$ ; (ii)  $\text{tr}\{(\mathbf{P}_s^* - \mathbf{P}_s)^\top(\mathbf{P}_s^* - \mathbf{P}_s)\} \leq C^2 p_s^2/n$ ; (iii)  $\text{tr}\{(\mathbf{P}_s^* - \mathbf{P}_s)^\top(\mathbf{P}_s^* - \mathbf{P}_s)\}^2 \leq Cp_s^4/n$ ; (iv)  $\mathbb{E}_{\max}(\mathbf{P}_s^*) \leq 1 + Cp_s/n$ ; (v)  $\text{tr}(\mathbf{P}_s^* \mathbf{P}_s^{*\top}) \leq Cp_s$ .*

**Proof.** We can obtain the proof using Lemma 3.1 of Ando and Li (2014).

**Lemma 3.** *(Hoeffding's inequality) Let  $X_1, \dots, X_n$  be independent random*



---

variables. Assume that  $\Pr(X_i \in [a_i, b_i]) = 1$  for  $1 \leq i \leq n$ , where  $a_i$  and  $b_i$  are constants. Let  $\bar{X} = \sum_{i=1}^n X_i/n$ . Then the following inequality holds

$$\Pr(|\bar{X} - E(\bar{X})| \geq t) \leq 2 \exp \left\{ - \frac{2n^2 t^2}{\sum_{i=1}^n (b_i - a_i)^2} \right\},$$

where  $t$  is a positive constant and  $E(\bar{X})$  is the expected value of  $X$ .

**Proof of Theorem 1.** Let  $C'$  be a constant. Denote  $L^*(\boldsymbol{\omega}) = \{\boldsymbol{\mu} - \tilde{\boldsymbol{\mu}}(\boldsymbol{\omega})\}^\top \widehat{\mathbf{W}} \{\boldsymbol{\mu} - \tilde{\boldsymbol{\mu}}(\boldsymbol{\omega})\}$  and  $\tilde{\boldsymbol{\mu}}(\boldsymbol{\omega}) = \tilde{\mathbf{P}}(\boldsymbol{\omega})\mathbf{Y}$ . Then, we have

$$\begin{aligned} \text{WCV}(\boldsymbol{\omega}) &= \{\mathbf{Y} - \tilde{\boldsymbol{\mu}}(\boldsymbol{\omega})\}^\top \widehat{\mathbf{W}} \{\mathbf{Y} - \tilde{\boldsymbol{\mu}}(\boldsymbol{\omega})\} \\ &= \{\boldsymbol{\mu} + \boldsymbol{\varepsilon} - \tilde{\mathbf{P}}(\boldsymbol{\omega})\mathbf{Y}\}^\top \widehat{\mathbf{W}} \{\boldsymbol{\mu} + \boldsymbol{\varepsilon} - \tilde{\mathbf{P}}(\boldsymbol{\omega})\mathbf{Y}\} \\ &= \boldsymbol{\varepsilon}^\top \widehat{\mathbf{W}} \boldsymbol{\varepsilon} + L^*(\boldsymbol{\omega}) + 2 \langle \widehat{\mathbf{W}}^{1/2} \boldsymbol{\varepsilon}, \widehat{\mathbf{W}}^{1/2} \{\boldsymbol{\mu} - \tilde{\mathbf{P}}(\boldsymbol{\omega})\mathbf{Y}\} \rangle \\ &= \boldsymbol{\varepsilon}^\top \widehat{\mathbf{W}} \boldsymbol{\varepsilon} + L(\boldsymbol{\omega}) \left\{ \frac{L^*(\boldsymbol{\omega})}{L(\boldsymbol{\omega})} + \frac{2 \langle \widehat{\mathbf{W}}^{1/2} \boldsymbol{\varepsilon}, \widehat{\mathbf{W}}^{1/2} \{\boldsymbol{\mu} - \tilde{\mathbf{P}}(\boldsymbol{\omega})\mathbf{Y}\} \rangle / R(\boldsymbol{\omega})}{L(\boldsymbol{\omega})/R(\boldsymbol{\omega})} \right\}. \end{aligned}$$

Thus,  $\hat{\boldsymbol{\omega}}$  can be obtained by minimizing  $\text{WCV}^*(\boldsymbol{\omega}) = \text{WCV}(\boldsymbol{\omega}) - \boldsymbol{\varepsilon}^\top \widehat{\mathbf{W}} \boldsymbol{\varepsilon}$  over  $\boldsymbol{\omega} \in \mathcal{W}$ . According to the definition of (2.4), if we can show

$$\sup_{\boldsymbol{\omega} \in \mathcal{W}} |L^*(\boldsymbol{\omega})/L(\boldsymbol{\omega}) - 1| \rightarrow 0, \quad (\text{S3.6})$$

$$\sup_{\boldsymbol{\omega} \in \mathcal{W}} | \langle \widehat{\mathbf{W}}^{1/2} \boldsymbol{\varepsilon}, \widehat{\mathbf{W}}^{1/2} \{\boldsymbol{\mu} - \tilde{\mathbf{P}}(\boldsymbol{\omega})\mathbf{Y}\} \rangle / R(\boldsymbol{\omega}) | \rightarrow 0, \quad (\text{S3.7})$$

$$\sup_{\boldsymbol{\omega} \in \mathcal{W}} |L(\boldsymbol{\omega})/R(\boldsymbol{\omega}) - 1| \rightarrow 0, \quad (\text{S3.8})$$

we can obtain that  $L(\hat{\boldsymbol{\omega}})/\inf_{\boldsymbol{\omega} \in \mathcal{W}} L(\boldsymbol{\omega}) \rightarrow 1$  is valid. Using the Cauchy-

Schwartz inequality leads to

$$\begin{aligned}
|L^*(\boldsymbol{\omega}) - L(\boldsymbol{\omega})| &= |\{\boldsymbol{\mu} - \tilde{\mathbf{P}}(\boldsymbol{\omega})\mathbf{Y}\}^\top \widehat{\mathbf{W}} \{\boldsymbol{\mu} - \tilde{\mathbf{P}}(\boldsymbol{\omega})\mathbf{Y}\} - \{\boldsymbol{\mu} - \widehat{\mathbf{P}}(\boldsymbol{\omega})\mathbf{Y}\}^\top \widehat{\mathbf{W}} \{\boldsymbol{\mu} - \widehat{\mathbf{P}}(\boldsymbol{\omega})\mathbf{Y}\}| \\
&= \|\widehat{\mathbf{W}}^{1/2} \{\tilde{\mathbf{P}}(\boldsymbol{\omega}) - \widehat{\mathbf{P}}(\boldsymbol{\omega})\}\mathbf{Y}\|^2 \\
&\quad - 2 \langle \widehat{\mathbf{W}}^{1/2} \{\boldsymbol{\mu} - \widehat{\mathbf{P}}(\boldsymbol{\omega})\mathbf{Y}\}, \widehat{\mathbf{W}}^{1/2} \{\tilde{\mathbf{P}}(\boldsymbol{\omega}) - \widehat{\mathbf{P}}(\boldsymbol{\omega})\}\mathbf{Y} \rangle | \\
&\leq \|\widehat{\mathbf{W}}^{1/2} \{\tilde{\mathbf{P}}(\boldsymbol{\omega}) - \widehat{\mathbf{P}}(\boldsymbol{\omega})\}\mathbf{Y}\|^2 + 2\sqrt{L(\boldsymbol{\omega})} \|\widehat{\mathbf{W}}^{1/2} \{\tilde{\mathbf{P}}(\boldsymbol{\omega}) - \widehat{\mathbf{P}}(\boldsymbol{\omega})\}\mathbf{Y}\|.
\end{aligned}$$

To show (S3.6), it is sufficient to show

$$\sup_{\boldsymbol{\omega} \in \mathcal{W}} \|\widehat{\mathbf{W}}^{1/2} \{\tilde{\mathbf{P}}(\boldsymbol{\omega}) - \widehat{\mathbf{P}}(\boldsymbol{\omega})\}\mathbf{Y}\|^2 / L(\boldsymbol{\omega}) \rightarrow 0. \quad (\text{S3.9})$$

From Lemma 1(iv), we have  $\tilde{\mathbf{P}}(\boldsymbol{\omega}) - \widehat{\mathbf{P}}(\boldsymbol{\omega}) = \mathbf{P}^*(\boldsymbol{\omega}) - \mathbf{P}(\boldsymbol{\omega}) + O_p(Sp_m^2 \sqrt{d_m/n})$ .

By (S3.8) and triangle inequality, the proof of (S3.9) is equivalent to proving

$$\sup_{\boldsymbol{\omega} \in \mathcal{W}} S^2 p_m^4 d_m n^{-1} / R(\boldsymbol{\omega}) \rightarrow 0, \quad (\text{S3.10})$$

$$\sup_{\boldsymbol{\omega} \in \mathcal{W}} \|\widehat{\mathbf{W}}^{1/2} \{\mathbf{P}^*(\boldsymbol{\omega}) - \mathbf{P}(\boldsymbol{\omega})\}\mathbf{Y}\|^2 / R(\boldsymbol{\omega}) \rightarrow 0. \quad (\text{S3.11})$$

Under Assumptions 3(iii), (iv) and (vi), we obtain

$$\begin{aligned}
\sup_{\boldsymbol{\omega} \in \mathcal{W}} S^2 p_m^4 d_m n^{-1} / R(\boldsymbol{\omega}) &\leq S^2 p_m^4 d_m n^{-1} / \xi_n \\
&= \left( \frac{S^{4G} \|\boldsymbol{\mu}\|^{2G}}{\xi_n^{2G}} \right)^{1/2G} \cdot \frac{\sqrt{n}}{\|\boldsymbol{\mu}\|} \cdot \left( \frac{p_m^{8/3} d_m}{n} \right)^{3/2} \cdot \frac{1}{\sqrt{d_m}} \rightarrow 0,
\end{aligned}$$

which indicates that (S3.10) holds. Applying the triangle inequality to

(S3.11) yields

$$\begin{aligned}
&\|\widehat{\mathbf{W}}^{1/2} \{\mathbf{P}^*(\boldsymbol{\omega}) - \mathbf{P}(\boldsymbol{\omega})\}\mathbf{Y}\|^2 \\
&= \left\{ \sum_{s=1}^S \omega_s \|\widehat{\mathbf{W}}^{1/2} (\mathbf{P}_s^* - \mathbf{P}_s) \boldsymbol{\mu}\| + \sum_{s=1}^S \omega_s \|\widehat{\mathbf{W}}^{1/2} (\mathbf{P}_s^* - \mathbf{P}_s) \boldsymbol{\varepsilon}\| \right\}^2
\end{aligned}$$

---


$$\begin{aligned}
&\leq \left\{ \sum_{s=1}^S \|\widehat{\mathbf{W}}^{1/2}(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\mu}\| + \sum_{s=1}^S \|\widehat{\mathbf{W}}^{1/2}(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\varepsilon}\| \right\}^2 \\
&\leq S^2 \left\{ \max_s \|\widehat{\mathbf{W}}^{1/2}(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\mu}\| + \max_s \|\widehat{\mathbf{W}}^{1/2}(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\varepsilon}\| \right\}^2 \\
&\leq 2S^2 \left\{ \max_s \|\widehat{\mathbf{W}}^{1/2}(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\mu}\|^2 + \max_s \|\widehat{\mathbf{W}}^{1/2}(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\varepsilon}\|^2 \right\}.
\end{aligned}$$

To prove (S3.11), it suffices to show that as  $n \rightarrow \infty$ , we have

$$S^2 \max_s \|\widehat{\mathbf{W}}^{1/2}(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\mu}\|^2 / \xi_n \rightarrow 0, \quad (\text{S3.12})$$

$$S^2 \max_s \|\widehat{\mathbf{W}}^{1/2}(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\varepsilon}\|^2 / \xi_n \rightarrow 0. \quad (\text{S3.13})$$

Note that

$$\frac{\delta_i}{\hat{\pi}_i} = \frac{\delta_i}{\pi_i} \left\{ 1 - \frac{\partial \gamma_1^\top \pi_i}{\pi_i} (\hat{\gamma}_1 - \gamma_{10}) + o_p(\sqrt{d_m/n}) \right\}.$$

By Theorem S1.1 and Assumption 1, we have

$$\sup_i \left| \frac{\delta_i}{\hat{\pi}_i} \right| \leq \sup_i \left| \frac{\delta_i}{\pi_i} \right| + \sup_i \left| \frac{\partial \gamma_1^\top \pi_i}{\pi_i^2} \right| \cdot \|\hat{\gamma}_1 - \gamma_{10}\|_2 \leq \frac{1}{C_0} + O_p(\sqrt{d_m/n}) \leq C'. \quad (\text{S3.14})$$

Therefore, for (S3.12) and (S3.13), it is sufficient to show

$$\begin{aligned}
S^2 \max_s \|\widehat{\mathbf{W}}^{1/2}(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\mu}\|^2 / \xi_n &\leq S^2 \max_s \sup_i |\delta_i / \hat{\pi}_i| \cdot \|(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\mu}\|^2 / \xi_n \\
&\leq C' S^2 \max_s \|(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\mu}\|^2 / \xi_n \rightarrow 0,
\end{aligned} \quad (\text{S3.15})$$

---


$$\begin{aligned}
S^2 \max_s \|\widehat{\mathbf{W}}^{1/2}(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\varepsilon}\|^2/\xi_n &\leq S^2 \max_s \sup_i |\delta_i/\widehat{\pi}_i| \cdot \|(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\varepsilon}\|^2/\xi_n \\
&\leq C' S^2 \max_s \|(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\varepsilon}\|^2/\xi_n \rightarrow 0.
\end{aligned} \tag{S3.16}$$

By Lemma 2(i), Assumptions 1 and 3(iii), (iv) and (vi), for any  $\kappa > 0$ , we only require showing

$$\begin{aligned}
&\Pr \left( S^2 \max_s \|(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\mu}\|^2/\xi_n > \kappa \right) \\
&\leq \sum_{s=1}^S \Pr \left( S^{4G} \|(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\mu}\|^{4G}/\xi_n^{2G} > \kappa^{2G} \right) \\
&\leq \frac{S^{4G}}{\kappa^{2G}\xi_n^{2G}} \sum_{s=1}^S \{E\|(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\mu}\|^{4G}\} \\
&\leq C' \frac{S^{4G}}{\kappa^{2G}\xi_n^{2G}} \sum_{s=1}^S \|(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\mu}\|^{4G} \\
&\leq C' \frac{S^{4G}}{\kappa^{2G}\xi_n^{2G}} \sum_{s=1}^S \{\mathbb{E}_{\max}(\mathbf{P}_s^* - \mathbf{P}_s)\}^{4G} \|\boldsymbol{\mu}\|^{4G} \\
&\leq C' C^{4G} \cdot \frac{S^{4G+1} \|\boldsymbol{\mu}\|^{2G}}{\xi_n^{2G} \kappa^{2G}} \cdot \left( \frac{\|\boldsymbol{\mu}\|^2}{n} \right)^G \cdot \left( \frac{p_m^{4/3}}{n} \right)^{3G} \rightarrow 0,
\end{aligned}$$

which implies that (S3.15) holds.

To prove (S3.16), it is sufficient to show that for any  $\kappa > 0$ , we have

$$\sum_{s=1}^S \Pr \left\{ S^2 \left| \|(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\varepsilon}\|^2 - E\|(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\varepsilon}\|^2 \right|/\xi_n > \kappa \right\} \rightarrow 0, \tag{S3.17}$$

$$S^2 \max_s E\|(\mathbf{P}_s^* - \mathbf{P}_s)\boldsymbol{\varepsilon}\|^2/\xi_n \rightarrow 0. \tag{S3.18}$$

By Theorem 2 of Whittle (1960) and Lemma 1(iii), it follows from (S3.17)

and Assumptions 3(i), (iii) and (vi) that

$$\begin{aligned}
& \sum_{s=1}^S \Pr \left\{ S^2 \left| \left\| (\mathbf{P}_s^* - \mathbf{P}_s) \boldsymbol{\varepsilon} \right\|^2 - E \left\| (\mathbf{P}_s^* - \mathbf{P}_s) \boldsymbol{\varepsilon} \right\|^2 \right| / \xi_n > \kappa \right\} \\
& \leq \sum_{s=1}^S S^{4G} E \left\{ \left\| (\mathbf{P}_s^* - \mathbf{P}_s) \boldsymbol{\varepsilon} \right\|^2 - E \left\| (\mathbf{P}_s^* - \mathbf{P}_s) \boldsymbol{\varepsilon} \right\|^2 \right\}^{2G} / (\xi_n^{2G} \kappa^{2G}) \\
& \leq C' \sum_{s=1}^S S^{4G} \left[ \text{tr} \{ (\mathbf{P}_s^* - \mathbf{P}_s)^\top (\mathbf{P}_s^* - \mathbf{P}_s) \} \right]^G / (\xi_n^{2G} \kappa^{2G}) \\
& \leq C' C^G \cdot \frac{S^{4G}}{\xi_n^{2G} \kappa^{2G}} \sum_{s=1}^S \left( \frac{p_s^4}{n^3} \right)^G \\
& \leq C' C^G \cdot \frac{S^{4G+1}}{\xi_n^{2G} \kappa^{2G}} \cdot \left( \frac{p_m^{4/3}}{n} \right)^{3G} \rightarrow 0,
\end{aligned}$$

which implies that (S3.17) holds. Also, it follows from Assumptions 3(iii), (vi) and Lemma 2(ii) that

$$\begin{aligned}
S^2 \max_s E \left\| (\mathbf{P}_s^* - \mathbf{P}_s) \boldsymbol{\varepsilon} \right\|^2 / \xi_n &= S^2 \max_s \text{tr} [ (\mathbf{P}_s^* - \mathbf{P}_s) \sigma_{\boldsymbol{\varepsilon}} (\mathbf{P}_s^* - \mathbf{P}_s)^\top ] \xi_n^{-1} \\
&\leq C^2 \mathbb{E}_{\max}(\sigma_{\boldsymbol{\varepsilon}}) \frac{p_s^2 S^2}{n \xi_n} \\
&= C^2 \cdot \mathbb{E}_{\max}(\sigma_{\boldsymbol{\varepsilon}}) \cdot \frac{p_s^2}{n} \cdot \left( \frac{S^{4G}}{\xi_n^{2G}} \right)^{1/2G} \rightarrow 0,
\end{aligned}$$

which implies that (S3.18) holds. Combining the proof of Lemma 1(iv) and triangle inequality, we obtain the following decomposition of (S3.7):

$$\begin{aligned}
& \left| \langle \widehat{\mathbf{W}}^{1/2} \boldsymbol{\varepsilon}, \widehat{\mathbf{W}}^{1/2} \{ \boldsymbol{\mu} - \widetilde{\mathbf{P}}(\boldsymbol{\omega}) \mathbf{Y} \} \rangle \right| \\
& \leq \left| \langle \widehat{\mathbf{W}}^{1/2} \boldsymbol{\varepsilon}, \widehat{\mathbf{W}}^{1/2} \boldsymbol{\mu} \rangle \right| + \left| \langle \widehat{\mathbf{W}}^{1/2} \boldsymbol{\varepsilon}, \widehat{\mathbf{W}}^{1/2} \mathbf{P}^*(\boldsymbol{\omega}) \boldsymbol{\mu} \rangle \right| \\
& \quad + \left| \langle \widehat{\mathbf{W}}^{1/2} \boldsymbol{\varepsilon}, \widehat{\mathbf{W}}^{1/2} \mathbf{P}^*(\boldsymbol{\omega}) \boldsymbol{\varepsilon} \rangle \right| + O_p \left( S p_m^2 \sqrt{\frac{d_m}{n}} \right).
\end{aligned}$$

---

Similar to the proof of (S3.10), under Assumptions 3(iii), (iv) and (vi), we delete the term  $O_p(Sp_m^2\sqrt{d_m/n})$ . Thus, the proof of (S3.7) is equivalent to proving

$$\sup_{\boldsymbol{\omega} \in \mathcal{W}} | \langle \widehat{\mathbf{W}}^{1/2} \boldsymbol{\varepsilon}, \widehat{\mathbf{W}}^{1/2} \boldsymbol{\mu} \rangle | / R(\boldsymbol{\omega}) \rightarrow 0, \quad (\text{S3.19})$$

$$\sup_{\boldsymbol{\omega} \in \mathcal{W}} | \langle \widehat{\mathbf{W}}^{1/2} \boldsymbol{\varepsilon}, \widehat{\mathbf{W}}^{1/2} \mathbf{P}^*(\boldsymbol{\omega}) \boldsymbol{\mu} \rangle | / R(\boldsymbol{\omega}) \rightarrow 0, \quad (\text{S3.20})$$

$$\sup_{\boldsymbol{\omega} \in \mathcal{W}} | \langle \widehat{\mathbf{W}}^{1/2} \boldsymbol{\varepsilon}, \widehat{\mathbf{W}}^{1/2} \mathbf{P}^*(\boldsymbol{\omega}) \boldsymbol{\varepsilon} \rangle | / R(\boldsymbol{\omega}) \rightarrow 0. \quad (\text{S3.21})$$

Using the similar arguments of (S3.14) and (S3.15), by Markov's inequality, under Assumptions 3(i), (vi) and (S3.14), given any  $\kappa > 0$ , we have

$$\begin{aligned} & \Pr \left\{ \sup_{\boldsymbol{\omega} \in \mathcal{W}} | \langle \widehat{\mathbf{W}}^{1/2} \boldsymbol{\varepsilon}, \widehat{\mathbf{W}}^{1/2} \boldsymbol{\mu} \rangle | / R(\boldsymbol{\omega}) > \kappa \right\} \\ & \leq \frac{E | \boldsymbol{\varepsilon}^\top \widehat{\mathbf{W}} \boldsymbol{\mu} |^{2G}}{\kappa^{2G} \xi_n^{2G}} \leq C' \kappa^{-2G} \frac{\| \boldsymbol{\mu} \|^{2G}}{\xi_n^{2G}} \rightarrow 0. \end{aligned}$$

For (S3.19), under (S3.14), we obtain  $| \boldsymbol{\varepsilon}^\top \widehat{\mathbf{W}} \mathbf{P}^*(\boldsymbol{\omega}) \boldsymbol{\mu} | \leq C' | \boldsymbol{\varepsilon}^\top \mathbf{P}^*(\boldsymbol{\omega}) \boldsymbol{\mu} |$ . According to Lemma 2(iv) and Assumption 3(i) and (vi), we just require proving

$$\begin{aligned} & \Pr \left( \sup_{\boldsymbol{\omega} \in \mathcal{W}} | \langle \boldsymbol{\varepsilon}, \mathbf{P}^*(\boldsymbol{\omega}) \boldsymbol{\mu} \rangle | / R(\boldsymbol{\omega}) > \kappa \right) \\ & \leq \Pr \left( S \max_s | \langle \boldsymbol{\varepsilon}, \mathbf{P}_s^* \boldsymbol{\mu} \rangle | > \kappa \xi_n \right) \\ & \leq \sum_{s=1}^S \Pr \left( S^{2G} | \langle \boldsymbol{\varepsilon}, \mathbf{P}_s^* \boldsymbol{\mu} \rangle |^{2G} > \kappa^{2G} \xi_n^{2G} \right) \end{aligned}$$

---


$$\begin{aligned}
&\leq \frac{S^{2G}}{\xi_n^{2G} \kappa^{2G}} \sum_{s=1}^S E |\boldsymbol{\varepsilon} \mathbf{P}_s^* \boldsymbol{\mu}|^{2G} \\
&\leq C' \frac{S^{2G}}{\xi_n^{2G} \kappa^{2G}} \sum_{s=1}^S \mathbb{E}_{\max}(\mathbf{P}_s^*)^{2G} \|\boldsymbol{\mu}\|^{2G} \\
&\leq C' \cdot (1 + C)^{2G} \cdot \frac{S^{2G+1} \|\boldsymbol{\mu}\|^{2G}}{\kappa^{2G} \xi_n^{2G}} \rightarrow 0.
\end{aligned}$$

Similarly, by Lemma 2(v), we only require proving

$$\begin{aligned}
&\Pr \left( \sup_{\boldsymbol{\omega} \in \mathcal{W}} | \langle \boldsymbol{\varepsilon}, \mathbf{P}^*(\boldsymbol{\omega}) \boldsymbol{\varepsilon} \rangle | / R(\boldsymbol{\omega}) > \kappa \right) \\
&\leq \Pr \left( S \max_s | \langle \boldsymbol{\varepsilon}, \mathbf{P}_s^* \boldsymbol{\varepsilon} \rangle | > \kappa \xi_n \right) \\
&\leq \frac{S^{2G}}{\kappa^{2G} \xi_n^{2G}} \sum_{s=1}^S E |\boldsymbol{\varepsilon}^\top \mathbf{P}_s^* \boldsymbol{\varepsilon}|^{2G} \\
&\leq C' \sigma_{\boldsymbol{\varepsilon}}^{2G} \frac{S^{2G}}{\kappa^{2G} \xi_n^{2G}} \sum_{s=1}^S \{\text{tr}(\mathbf{P}_s^* \mathbf{P}_s^{*\top})\}^G \\
&\leq C' \cdot C^G \cdot \sigma_{\boldsymbol{\varepsilon}}^{2G} \cdot \frac{S^{2G+1} n^G}{\kappa^{2G} \xi_n^{2G}} \rightarrow 0,
\end{aligned}$$

where the third inequality holds because of Assumption 1 and Assumption 3(i), the last term converges to zero because of Assumption 3(vi).

To show (S3.8), using the similar arguments of (S3.14) and (S3.15), by

Triangle inequality, we have

$$\begin{aligned}
&|L(\boldsymbol{\omega}) - R(\boldsymbol{\omega})| \\
&= | \{\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}(\boldsymbol{\omega})\}^\top \widehat{\mathbf{W}} \{\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}(\boldsymbol{\omega})\} - E[\{\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}(\boldsymbol{\omega})\}^\top \widehat{\mathbf{W}} \{\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}(\boldsymbol{\omega})\} | \mathbf{X}] | \\
&\leq \| \widehat{\mathbf{W}}^{1/2} \widetilde{\mathbf{M}}(\boldsymbol{\omega}) \boldsymbol{\mu} \|^2 + \| \widehat{\mathbf{W}}^{1/2} \widehat{\mathbf{P}}(\boldsymbol{\omega}) \boldsymbol{\varepsilon} \|^2 + 2 | \langle \widehat{\mathbf{W}}^{1/2} \widetilde{\mathbf{M}}(\boldsymbol{\omega}) \boldsymbol{\mu}, \widehat{\mathbf{W}}^{1/2} \widehat{\mathbf{P}}(\boldsymbol{\omega}) \boldsymbol{\varepsilon} \rangle |
\end{aligned}$$

---


$$\begin{aligned}
& + E\{\|\widehat{\mathbf{W}}^{1/2}\widetilde{\mathbf{M}}(\boldsymbol{\omega})\boldsymbol{\mu}\|^2|\mathbf{X}\} + E\{\|\widehat{\mathbf{W}}^{1/2}\widehat{\mathbf{P}}(\boldsymbol{\omega})\boldsymbol{\varepsilon}\|^2|\mathbf{X}\} \\
& + |E\{2\langle \widehat{\mathbf{W}}^{1/2}\widetilde{\mathbf{M}}(\boldsymbol{\omega})\boldsymbol{\mu}, \widehat{\mathbf{W}}^{1/2}\widehat{\mathbf{P}}(\boldsymbol{\omega})\boldsymbol{\varepsilon} \rangle|\mathbf{X}\}| \\
& \leq C' \|\widetilde{\mathbf{M}}(\boldsymbol{\omega})\boldsymbol{\mu}\|^2 + C' \|\widehat{\mathbf{P}}(\boldsymbol{\omega})\boldsymbol{\varepsilon}\|^2 + 2C' |\langle \widetilde{\mathbf{M}}(\boldsymbol{\omega})\boldsymbol{\mu}, \widehat{\mathbf{P}}(\boldsymbol{\omega})\boldsymbol{\varepsilon} \rangle| \\
& + C' E\{\|\widetilde{\mathbf{M}}(\boldsymbol{\omega})\boldsymbol{\mu}\|^2|\mathbf{X}\} + C' E\{\|\widehat{\mathbf{P}}(\boldsymbol{\omega})\boldsymbol{\varepsilon}\|^2|\mathbf{X}\} + 2C' |E\{\langle \widetilde{\mathbf{M}}(\boldsymbol{\omega})\boldsymbol{\mu}, \widehat{\mathbf{P}}(\boldsymbol{\omega})\boldsymbol{\varepsilon} \rangle|\mathbf{X}\}|
\end{aligned}$$

where  $\widetilde{\mathbf{M}}(\boldsymbol{\omega}) = \mathbf{I} - \widehat{\mathbf{P}}(\boldsymbol{\omega})$ . From the proof of Lemma 1(iv), by (S3.10) and Assumption 3(vi), if we delete the remainder term  $O_p(S^2 p_m^4 d_m/n)$ , the proof of (S3.8) is equivalent to show,

$$\sup_{\boldsymbol{\omega} \in \mathcal{W}} \left| \frac{\|\mathbf{M}(\boldsymbol{\omega})\boldsymbol{\mu}\|^2}{R(\boldsymbol{\omega})} \right| \rightarrow 0, \quad (\text{S3.22})$$

$$\sup_{\boldsymbol{\omega} \in \mathcal{W}} \left| \frac{\|\mathbf{P}(\boldsymbol{\omega})\boldsymbol{\varepsilon}\|^2 - \text{tr}\{\mathbf{P}(\boldsymbol{\omega})\boldsymbol{\sigma}_{\boldsymbol{\varepsilon}}\mathbf{P}^\top(\boldsymbol{\omega})\}}{R(\boldsymbol{\omega})} \right| \rightarrow 0, \quad (\text{S3.23})$$

$$\sup_{\boldsymbol{\omega} \in \mathcal{W}} \left| \frac{\langle \mathbf{M}(\boldsymbol{\omega})\boldsymbol{\mu}, \mathbf{P}(\boldsymbol{\omega})\boldsymbol{\varepsilon} \rangle}{R(\boldsymbol{\omega})} \right| \rightarrow 0, \quad (\text{S3.24})$$

where  $\mathbf{M} = \mathbf{I} - \mathbf{P}(\boldsymbol{\omega})$ . Based on Assumption 3(vi), we can obtain (S3.22) using the method given in the proof of (S3.12). Next, we prove (S3.23) under Assumptions 3(i), (iv) and (vi). For any  $\kappa > 0$ , we have

$$\begin{aligned}
& \Pr \left\{ \sup_{\boldsymbol{\omega} \in \mathcal{W}} \left| \frac{\|\mathbf{P}(\boldsymbol{\omega})\boldsymbol{\varepsilon}\|^2 - \text{tr}\{\mathbf{P}(\boldsymbol{\omega})\boldsymbol{\sigma}_{\boldsymbol{\varepsilon}}\mathbf{P}^\top(\boldsymbol{\omega})\}}{R(\boldsymbol{\omega})} \right| > \kappa \right\} \\
& \leq \Pr \left\{ \sup_{\boldsymbol{\omega} \in \mathcal{W}} \sum_{s=1}^S \sum_{k=1}^S \omega_s \omega_k |\boldsymbol{\varepsilon}^\top \mathbf{P}_s \mathbf{P}_k \boldsymbol{\varepsilon} - \boldsymbol{\sigma}_{\boldsymbol{\varepsilon}} \text{tr}(\mathbf{P}_s \mathbf{P}_k)| > \kappa \xi_n \right\} \\
& \leq \Pr \left\{ S^2 \max_s \max_k |\boldsymbol{\varepsilon}^\top \mathbf{P}_s \mathbf{P}_k \boldsymbol{\varepsilon} - \boldsymbol{\sigma}_{\boldsymbol{\varepsilon}} \text{tr}(\mathbf{P}_s \mathbf{P}_k)| > \kappa \xi_n \right\}
\end{aligned}$$



---


$$\begin{aligned}
&\leq \sum_{s=1}^S \sum_{k=1}^S \Pr \left\{ \left| \boldsymbol{\varepsilon}^\top \mathbf{P}_s \mathbf{P}_k \boldsymbol{\varepsilon} - \sigma_{\boldsymbol{\varepsilon}} \text{tr} \mathbf{P}_s \mathbf{P}_k \right| > \kappa \xi_n / S^2 \right\} \\
&\leq \sum_{s=1}^S \sum_{k=1}^S \frac{S^{4G}}{\kappa^{2G} \xi_n^{2G}} E \left| \boldsymbol{\varepsilon}^\top \mathbf{P}_s \mathbf{P}_k \boldsymbol{\varepsilon} - \sigma_{\boldsymbol{\varepsilon}} \text{tr} \mathbf{P}_s \mathbf{P}_k \right|^{2G} \\
&\leq C' \frac{S^{4G}}{\kappa^{2G} \xi_n^{2G}} \sum_{s=1}^S \sum_{k=1}^S \left\{ \text{tr}(\mathbf{P}_s^2 \mathbf{P}_k^2) \right\}^G \\
&\leq C' \cdot \frac{S^{4G+2} n^G}{\xi_n^{2G}} \cdot \kappa^{-2G} \rightarrow 0,
\end{aligned}$$

where the last inequality holds since  $\mathbf{P}(\boldsymbol{\omega})$  is the idempotent matrix. Similarly,  $\mathbf{I} - \mathbf{P}(\boldsymbol{\omega})$  is also the idempotent matrix. Then, under Assumptions 3(i) and (vi), for any  $\kappa > 0$ , we have

$$\begin{aligned}
&\Pr \left\{ \sup_{\boldsymbol{\omega} \in \mathcal{W}} \left| \frac{\langle \mathbf{M}(\boldsymbol{\omega}) \boldsymbol{\mu}, \mathbf{P}(\boldsymbol{\omega}) \boldsymbol{\varepsilon} \rangle}{R(\boldsymbol{\omega})} \right| > \kappa \right\} \\
&\leq \Pr \left\{ \sup_{\boldsymbol{\omega} \in \mathcal{W}} \sum_{s=1}^S \sum_{k=1}^S \left| \boldsymbol{\mu}^\top (\mathbf{I} - \mathbf{P}_s) \mathbf{P}_k \boldsymbol{\varepsilon} \right| > \kappa \xi_n \right\} \\
&\leq \frac{S^{4G}}{\kappa^{2G} \xi_n^{2G}} \sum_{s=1}^S \sum_{k=1}^S E \left| \boldsymbol{\mu}^\top (\mathbf{I} - \mathbf{P}_s) \mathbf{P}_k \boldsymbol{\varepsilon} \right|^{2G} \\
&\leq C' \frac{S^{4G}}{\kappa^{2G} \xi_n^{2G}} \sum_{s=1}^S \sum_{k=1}^S \left\| \mathbf{P}_s (\mathbf{I} - \mathbf{P}_k) \boldsymbol{\mu} \right\|^{2G} \\
&\leq C' \cdot \frac{S^{4G+2} \|\boldsymbol{\mu}\|^{2G}}{\xi_n^{2G}} \cdot \kappa^{-2G} \rightarrow 0,
\end{aligned}$$

which indicates that we have proved (S3.21). Thus, we finish the proof of Theorem 1.

**Proof of Theorem S2.1.** We first show the sure screening property. For

each  $k = 1, \dots, p$ , define

$$m(x) = E(X_k Y | X_k = x)$$

and

$$\hat{m}(x) = \frac{\sum_{j=1}^n \delta_j K_h(x - X_{jk}) Y_j}{\sum_{j=1}^n \delta_j K_h(x - X_{jk})}.$$

By the definition of  $\hat{r}_k$  and  $r_k$ , we have

$$\begin{aligned} \hat{r}_k - r_k &= \frac{1}{n} \sum_{i=1}^n \left\{ \delta_i X_{ik} Y_i + (1 - \delta_i) \frac{1}{m} \sum_{v=1}^m X_{ik} \tilde{Y}_{iv}^k \right\} - E(X_k Y) \\ &= \frac{1}{n} \sum_{i=1}^n (1 - \delta_i) \{ \hat{m}(X_{ik}) - m(X_{ik}) \} + \frac{1}{n} \sum_{i=1}^n (1 - \delta_i) \left\{ \frac{1}{m} \sum_{v=1}^m X_{ik} \tilde{Y}_{iv}^k - \hat{m}(X_{ik}) \right\} \\ &\quad + \frac{1}{n} \sum_{i=1}^n \delta_i \{ X_{ik} Y_i - m(X_{ik}) \} + \frac{1}{n} \sum_{i=1}^n \{ m(X_{ik}) - E(X_k Y) \} \\ &= J_{k1} + J_{k2} + J_{k3} + J_{k4}. \end{aligned}$$

For  $J_{k2}$ , Wang and Chen (2009) have proved  $J_{k3} = o_p(1/\sqrt{n})$  as  $n$  and  $m \rightarrow \infty$ . Thus, as  $n$  is sufficiently large, for any  $\varsigma \in (0, 1/2)$  and  $c_2 > 0$ , we have

$$\begin{aligned} \Pr\{|\hat{r}_k - r_k| \geq c_2 n^{-\varsigma}\} &= \Pr(|J_{k1} + J_{k2} + J_{k3} + J_{k4}| \geq c_2 n^{-\varsigma}) \\ &\leq \Pr(|J_{k1} + J_{k3} + J_{k4}| \geq c_2 n^{-\varsigma} - |J_{k2}|) \\ &\leq \Pr(|J_{k1} + J_{k3} + J_{k4}| \geq c_2 n^{-\varsigma}/2). \end{aligned}$$

Through defining  $\eta(x) = \pi(x) f_k(x)$  and  $\hat{\eta}(x) = \sum_{j=1}^n \delta_j K_h(X_{jk} - x)/n$ , where  $\pi(\cdot)$  is the selection probability function and  $f_k(\cdot)$  is the probability density function of  $X_k$ . Thus we can decompose  $J_{k1}$  as

---


$$\begin{aligned}
J_{k1} &= \frac{1}{n} \sum_{i=1}^n (1 - \delta_i) \{\hat{m}(X_{ik}) - m(X_{ik})\} \\
&= \frac{1}{n} \sum_{i=1}^n (1 - \delta_i) \frac{\sum_{j=1}^n \delta_j K_h(X_{jk} - X_{ik}) \{X_{ik} Y_i - m(X_{jk})\} / n}{\eta(X_{ik})} \\
&\quad + \frac{1}{n} \sum_{i=1}^n (1 - \delta_i) \{\hat{m}(X_{ik}) - m(X_{ik})\} \frac{\eta(X_{ik}) - \hat{\eta}(X_{ik})}{\eta(X_{ik})} \\
&\quad + \frac{1}{n} \sum_{i=1}^n (1 - \delta_i) \frac{\sum_{j=1}^n \delta_j K_h(X_{jk} - X_{ik}) \{m(X_{jk}) - m(X_{ik})\} / n}{\eta(X_{ik})} \\
&= J_{k11} + J_{k12} + J_{k13} \\
&= \tilde{J}_{k11} + (J_{k11} - \tilde{J}_{k11}) + J_{k12} + J_{k13},
\end{aligned}$$

where  $\tilde{J}_{k11} = (1/n) \sum_{i=1}^n \delta_i \{X_{ik} Y_i - m(X_{ik})\} \{1 - \pi(X_{ik})\} / \pi(X_{ik})$ . Under the Assumptions 1, 8, and 9, by the similar certification of Wang and Chen (2009),  $J_{k11} - \tilde{J}_{k11} = o_p(1/\sqrt{n})$ ,  $J_{k12} = o_p(1/\sqrt{n})$  and  $J_{k13} = o_p(1/\sqrt{n})$ .

Then we can write

$$\begin{aligned}
&\Pr(|J_{k1} + J_{k3} + J_{k4}| \geq c_2 n^{-\varsigma} / 2) \\
&\leq \Pr(|\tilde{J}_{k11} + (J_{k11} - \tilde{J}_{k11}) + J_{k12} + J_{k13} + J_{k3} + J_{k4}| \geq c_2 n^{-\varsigma} / 2) \\
&\leq \Pr(|\tilde{J}_{k11} + J_{k3} + J_{k4}| \geq c_2 n^{-\varsigma} / 2 - |J_{k11} - \tilde{J}_{k11}| - |J_{k12}| - |J_{k13}|) \\
&\leq \Pr(|\tilde{J}_{k11} + J_{k3} + J_{k4}| \geq c_2 n^{-\varsigma} / 16).
\end{aligned}$$

---

Note that

$$\begin{aligned}
& \tilde{J}_{k11} + J_{k3} + J_{k4} \\
&= (1/n) \sum_{i=1}^n \delta_i \{X_{ik} Y_i - m(X_{ik})\} \{1 - \pi(X_{ik})\} / \pi(X_{ik}) \\
&\quad + \frac{1}{n} \sum_{i=1}^n \delta_i \{X_{ik} Y_i - m(X_{ik})\} + \frac{1}{n} \sum_{i=1}^n \{m(X_{ik}) - E(X_k Y)\} \\
&= \frac{1}{n} \sum_{i=1}^n \delta_i \{X_{ik} Y_i - E(X_k Y)\} \frac{1}{\pi(X_{ik})} \\
&\quad + \frac{1}{n} \sum_{i=1}^n \frac{\{\pi(X_{ik}) - \delta_i\}}{\pi(X_{ik})} [m(X_{ik}) - E\{m(X_{ik})\}] \\
&= I_1 + I_2.
\end{aligned}$$

Under Assumption 1, for any  $M > 0$ , we have

$$\begin{aligned}
\Pr(|I_1| \geq c_2 n^{-\varsigma} / 32) &= \Pr \left\{ \left| \frac{1}{n} \sum_{i=1}^n \delta_i \{X_{ik} Y_i - E(X_k Y)\} \frac{1}{\pi(X_{ik})} \right| \geq c_2 n^{-\varsigma} / 32 \right\} \\
&\leq \Pr \left\{ \left| \frac{1}{n} \sum_{i=1}^n \delta_i \{X_{ik} Y_i - E(X_k Y)\} \right| \geq C_0 c_2 n^{-\varsigma} / 32 \right\} \\
&\leq \Pr \left\{ \left| \frac{1}{n} \sum_{i=1}^n \delta_i \{X_{ik} Y_i - E(X_k Y)\} \right| \geq C_0 c_2 n^{-\varsigma} / 32, \max_i |\delta_i X_{ik} Y_i| < M \right\} \\
&\quad + \Pr(\max_i |X_{ik} Y_i| \geq M).
\end{aligned}$$

According to Assumption 10 and Lemma S3 of Liu et al. (2014), there are some positive constants  $c_3$  and  $c_4$  such that for any  $M > 0$ , we have  $\Pr(|X_k Y| \geq M) \leq c_3 \exp(-c_4 M)$ . Thus, by taking  $M = c_3^{-1} n^{(1-2\varsigma)/3}$ , applying the Hoeffdings inequality in Lemma 3 and yields that there exists

---

a positive constant  $c_5$  such that

$$\begin{aligned}
& \Pr \left\{ \left| \frac{1}{n} \sum_{i=1}^n \delta_i \{X_{ik} Y_i - E(X_k Y)\} \right| \geq C_0 c_2 n^{-\varsigma} / 32, \max_i |\delta_i X_{ik} Y_i| < M \right\} \\
& + \Pr(\max_i |X_{ik} Y_i| \geq M) \\
& \leq 2 \exp(-n^{1-2\varsigma} / M^2) + \sum_{i=1}^n \Pr(|X_{ik} Y_i| \geq M) \\
& \leq O(n) \exp\{-c_5 n^{(1-2\varsigma)/3}\}.
\end{aligned}$$

According to the Assumptions 1 and 10, the above argument can be used to  $I_2$ , we have  $\Pr(|I_2| \geq c_2 n^{-\varsigma} / 32) \leq O(n) \exp\{-c_6 n^{(1-2\varsigma)/3}\}$  for some constant  $c_6$ . Thus, there exists a constant  $c_7$  such that

$$\begin{aligned}
& \Pr \left( \max_{1 \leq k \leq p} |\hat{r}_k - r_k| \geq c_2 n^{-\varsigma} \right) \\
& \leq p \Pr(|I_1| \geq c_2 n^{-\varsigma} / 32) + p \Pr(|I_2| \geq c_2 n^{-\varsigma} / 32) \\
& \leq O(n) p \exp\{-c_7 n^{(1-2\varsigma)/3}\} = O\{p \exp(-c_7 n^{(1-2\varsigma)/3} + \log(n))\}.
\end{aligned}$$

In fact, by  $S_n = \{\max_{k \in \mathcal{M}_*} |\hat{r}_k - r_k| \leq c_0 n^{-\varsigma} / 2\}$  and Assumption 11, we have

$$\min_{k \in \mathcal{M}_*} |\hat{r}_k| \geq \min_{k \in \mathcal{M}_*} (|r_k| - |\hat{r}_k - r_k|) \geq \min_{k \in \mathcal{M}_*} |r_k| - \max_{k \in \mathcal{M}_*} |\hat{r}_k - r_k| \geq c_0 n^{-\varsigma} / 2.$$

Thus, by taking  $\varrho_n = c_8 n^{-\varsigma}$  with  $c_8 \leq c_0 / 2$ , there exists some positive constant  $c_9$  such that

$$\Pr(\mathcal{M}_* \subset \widehat{\mathcal{M}}) \geq \Pr(S_n) \geq 1 - O\{|\mathcal{M}_*| \exp(-c_9 n^{(1-2\varsigma)/3} + \log(n))\}.$$

Now we show the ranking consistency property. It is easily shown that

$$\begin{aligned}
& \Pr \left\{ \left( \min_{k \in \mathcal{M}_*} |\hat{r}_k| - \max_{k \in \mathcal{I}_*} |\hat{r}_k| \right) < m_0/2 \right\} \\
& \leq \Pr \left\{ \left( \min_{k \in \mathcal{M}_*} |\hat{r}_k| - \max_{k \in \mathcal{I}_*} |\hat{r}_k| \right) - \left( \min_{k \in \mathcal{M}_*} |r_k| - \max_{k \in \mathcal{I}_*} |r_k| \right) < -m_0/2 \right\} \\
& \leq \Pr \left\{ \left| \left( \min_{k \in \mathcal{M}_*} |\hat{r}_k| - \max_{k \in \mathcal{I}_*} |\hat{r}_k| \right) - \left( \min_{k \in \mathcal{M}_*} |r_k| - \max_{k \in \mathcal{I}_*} |r_k| \right) \right| \geq m_0/2 \right\} \\
& \leq \Pr \left( 2 \max_{1 \leq k \leq p} |\hat{r}_k - r_k| \geq m_0/2 \right) \\
& \leq p \Pr(|\hat{r}_k - r_k| \geq m_0/4) \\
& \leq O(n)p \exp(-c_{10}n^{1/3}m_0^{2/3})
\end{aligned}$$

for some constants  $c_{10}$ , where the first inequality holds because of Assumption 12.

Note that  $\log(n) = o(n^{1/3}m_0^{2/3})$  and  $\log(p) = o(n^{1/3}m_0^{2/3})$  imply that  $p \leq \exp(c_{10}n^{1/3}m_0^{2/3}/2)$ ,  $c_{10}n^{1/3}m_0^{2/3}/2 \geq 4 \log(n)$  for sufficiently large  $n$ .

Thus, for some  $n_0$ , we have

$$\begin{aligned}
\sum_{n=n_0}^{\infty} pn \exp(-c_{10}n^{1/3}m_0^{2/3}) & \leq 2 \sum_{n=n_0}^{\infty} \exp\{c_{10}n^{1/3}m_0^{2/3}/2 - c_{10}n^{1/3}m_0^{2/3} + \log(n)\} \\
& \leq 2 \sum_{n=n_0}^{\infty} \exp\{-3 \log(n)\} \leq 2 \sum_{n=n_0}^{\infty} n^{-3} < +\infty.
\end{aligned}$$

Hence, by Borel Contelli Lemma, we have

$$\liminf_{n \rightarrow \infty} \left\{ \min_{k \in \mathcal{M}_*} |\hat{r}_k| - \max_{k \in \mathcal{I}_*} |\hat{r}_k| \right\} > 0 \quad a.s.$$

## References

Cui, H., Li, R. and Zhong, W. (2015). Model-free feature screening for ultrahigh dimensional discriminant analysis. *Journal of the American Statistical Association* **110**, 630-641.

- 
- Fan, J. and Lv, J. (2008). Sure independence screening for ultrahigh dimensional feature space. *Journal of the Royal Statistical Society, Series B* **70**, 849-911.
- Fan, J. and Lv, J. (2011). Nonconcave penalized likelihood with NP-dimensionality. *IEEE Transactions on Information Theory* **57**, 5467-5484.
- Fan, J. and Peng, H. (2004). Nonconcave penalized likelihood with a diverging number of parameters. *The Annals of Statistics* **32**, 928-961.
- Hirano, K., Imbens, G. and Ridder, G. (2003). Efficient estimation of average treatment effects using the estimated propensity score. *Econometrica* **71**, 1161-1189.
- Lewis, R. and Reinsel, G. C. (1985). Prediction of multivariate time series by autoregressive model fitting. *Journal of Multivariate Analysis* **16**, 393-411.
- Li, R., Zhong, W. and Zhu, L. (2012). Feature screening via distance correlation learning. *Journal of the American Statistical Association* **107**, 1129-1139.
- Liu, J., Li, R. and Wu, R. (2014). Feature selection for varying coefficient models with ultrahigh dimensional covariates. *Journal of the American Statistical Association* **109**, 266-274.
- Lv, J. and Fan, Y. (2009). A unified approach to model selection and sparse recovery using regularized least squares. *The Annals of Statistics* **37**, 3498-3528.
- Whittle, P. (1960). Bounds for the moments of linear and quadratic forms in independent variables. *Theory of Probability and its Applications* **5**, 302-305.
- Wiener, N. and Masani, P. (1958). The prediction theory of multivariate stochastic processes,

II. The linear predictor. *Acta Mathematica* **99**, 93-137.

Zhao, P. and Yu, B. (2006). On model selection consistency of lasso. *Journal of Machine Learning Research* **7**, 2541-2567.

Zhang, C. H. (2010). Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics* **38**, 894-942.