

**SUPPLEMENT TO “LINEAR HYPOTHESIS TESTING
FOR HIGH DIMENSIONAL GENERALIZED LINEAR
MODELS”**

BY CHENGCHUN SHI^{*}, RUI SONG^{*}, ZHAO CHEN[†], AND RUNZE LI[†]

North Carolina State University and Pennsylvania State University

In this Web Appendix, we compare the power of the proposed test statistics with existing tests in the literature in Section S1. Section S2 proves the convergence rate of $\hat{\phi}$ proposed in Section 3.3.2. Section S3 discusses the asymptotic powers of partial penalized test statistics. Discussions of Condition (A1), (A2), (A3) and (A4) are presented in Section S4. Proofs and technical lemmas are given in Section S5 and Section S6 respectively. Section S7 contains the real data analysis. Section S8 includes additional simulation studies for Poisson regression model and additional tables and plots.

APPENDIX S1: POWER COMPARISONS

In this section, we consider the following class of null hypothesis: $H_0 : \beta_{0,j} = 0$ and compare the power of our tests with the Wald test in [van de Geer et al. \(2014\)](#) and the decorrelated score test statistic in [Ning and Liu \(2017\)](#). Without loss of generality, we fix $j = 1$. Consider the following local alternative hypothesis $H_a : \beta_{0,1} = n^{-1/2}h$ for some $h \neq 0$. According to Corollary 3.1, with some calculations, we can show that for $T = T_W, T_S$ and T_L , the power function takes the form

$$\begin{aligned} \Pr(T > \chi_\alpha^2(1)) &= \Pr(\chi^2(r, \gamma_n) > \chi_\alpha^2(1)) + o(1) = \Pr(|\mathbb{Z} + \sqrt{\gamma_n}| > z_{\frac{\alpha}{2}}) + o(1) \\ &= \Pr(\mathbb{Z} > z_{\frac{\alpha}{2}} - \sqrt{\gamma_n}) + \Pr(\mathbb{Z} < -z_{\frac{\alpha}{2}} - \sqrt{\gamma_n}) + o(1) \\ \text{(S1.1)} \quad &= g(\alpha, \gamma_n) + o(1), \end{aligned}$$

where $g(\alpha, \gamma) = \Pr(\mathbb{Z} > z_{\alpha/2} - \sqrt{\gamma}) + \Pr(\mathbb{Z} \geq z_{\alpha/2} + \sqrt{\gamma})$ for any $\gamma \geq 0$, \mathbb{Z} denotes a standard normal variable, z_α is the upper α th quantile of a standard normal distribution, and $\gamma_n = (\phi_0 \mathbf{e}_{1,1+s}^T \boldsymbol{\Omega}_n \mathbf{e}_{1,1+s})^{-1} h^2$ where $\mathbf{e}_{1,q}$ denotes the basis vector of length q , with the first element equal to 1 and other elements equal to 0. To simplify the presentation, we consider a

^{*}Supported by NSF grant DMS 1555244, NCI grant P01 CA142538.

[†]Supported by NSF grant DMS 1512422, NIH grants P50 DA039838 and P50 DA036107, and T32 LM012415, and NNSFC grants 11690014 and 11690015.

random design setting where $\mathbf{X}_1, \dots, \mathbf{X}_n$ are independent and identically distributed, according to the distribution of \mathbf{X}_0 . Define

$$\Phi = \mathbb{E} \{ \mathbf{X}_0 b''(\mathbf{X}_0 \beta_0) \mathbf{X}_0^T \}.$$

Under certain regularity conditions, it follows from (S1.1) that

$$\Pr(T > \chi_\alpha^2(1)) = g(\alpha, \{\phi_0 \mathbf{e}_{1,1+s}^T (\Phi_{\{1\} \cup S, \{1\} \cup S})^{-1} \mathbf{e}_{1,1+s}\}^{-1} h^2) + o(1).$$

S1.0.1. *Comparison with the Wald test based on the de-sparsified Lasso estimator.* van de Geer et al. (2014) proposed to de-sparsify the Lasso estimator for constructing statistical tests for $\beta_{0,j}$ in a high dimensional generalized linear model. Define the Lasso estimator

$$\hat{\beta}^L = \arg \min_{\beta} \left(\frac{1}{n} \sum_{i=1}^n \{b(\beta^T \mathbf{X}_i) - Y_i \beta^T \mathbf{X}_i\} + \lambda_n \|\beta\|_1 \right).$$

The de-sparsified Lasso estimator $\hat{\beta}^{DL}$ is given by

$$\hat{\beta}^{DL} = \hat{\beta}^L + \frac{1}{n} \hat{\Theta} \mathbf{X}^T \{ \mathbf{Y} - \mu(\mathbf{X} \hat{\beta}^L) \},$$

where $\hat{\Theta}$ is some consistent estimator for Φ^{-1} obtained by nodewise regression. It follows from Theorem 3.1 in van de Geer et al. (2014) that

$$(S1.2) \quad \sqrt{n}(\hat{\beta}_j^{DL} - \beta_{0,j})/\hat{\sigma}_j \sim N(0, 1),$$

where

$$\hat{\sigma}_j^2 = \frac{1}{n} \left(\mathbf{e}_{j,p}^T \hat{\Theta} \mathbf{X}^T \{ \mathbf{Y} - \mu(\mathbf{X} \hat{\beta}^L) \} \{ \mathbf{Y} - \mu(\mathbf{X} \hat{\beta}^L) \}^T \mathbf{X} \hat{\Theta} \mathbf{e}_{j,p} \right).$$

Moreover, under the model assumption (1.1), it follows from Corollary 3.1 and the proof of Theorem 3.1 in van de Geer et al. (2014) that

$$(S1.3) \quad \hat{\sigma}_j^2 = \phi_0 \mathbf{e}_{j,p}^T \Phi^{-1} \mathbf{e}_{j,p} + o_p(1).$$

To test $H_0 : \beta_{0,1} = 0$, van de Geer et al. (2014) considered the Wald-type statistic and proposed to reject H_0 when $\sqrt{n}|\hat{\beta}_1^{DL}| > z_{\alpha/2} \hat{\sigma}_1$ for a given α . It follows from (S1.2) that such test statistic has correct size under H_0 . Assume $\liminf_n \mathbf{e}_{1,p}^T \Phi^{-1} \mathbf{e}_{1,p} > 0$, under the local alternative, it follows from (S1.2)

and (S1.3) that

$$\begin{aligned}
& \Pr(\sqrt{n}|\hat{\beta}_1^{DL}| > z_{\alpha/2}\hat{\sigma}_1) = \Pr(\sqrt{n}(\hat{\beta}_1^{DL} - \beta_{0,1}) > z_{\alpha/2}\hat{\sigma}_1 - h) \\
& + \Pr(\sqrt{n}(\hat{\beta}_1^{DL} - \beta_{0,1}) < -z_{\alpha/2}\hat{\sigma}_1 - h) \\
& = \Pr\left(\frac{\sqrt{n}(\hat{\beta}_1^{DL} - \beta_{0,1})}{\hat{\sigma}_1} > z_{\alpha/2} - \frac{h}{\sqrt{\phi_0 \mathbf{e}_{1,p}^T \mathbf{\Phi}^{-1} \mathbf{e}_{1,p}}}\right) \\
& + \Pr\left(\frac{\sqrt{n}(\hat{\beta}_1^{DL} - \beta_{0,1})}{\hat{\sigma}_1} < -z_{\alpha/2} - \frac{h}{\sqrt{\phi_0 \mathbf{e}_{1,p}^T \mathbf{\Phi}^{-1} \mathbf{e}_{1,p}}}\right) + o(1) \\
& = g(\alpha, (\phi_0 \mathbf{e}_{1,p}^T \mathbf{\Phi}^{-1} \mathbf{e}_{1,p})^{-1} h^2) + o(1).
\end{aligned}$$

LEMMA S.1. For any $p \times p$ positive definite matrix $\mathbf{\Phi}$ and any set $S \subseteq [2, 3, \dots, p]$, we have

$$\mathbf{e}_{1,p}^T \mathbf{\Phi}^{-1} \mathbf{e}_{1,p} \geq \mathbf{e}_{1,1+s}^T (\mathbf{\Phi}_{\{1\} \cup S, \{1\} \cup S})^{-1} \mathbf{e}_{1,1+s},$$

with equality if and only if

$$\mathbf{\Phi}_{S, \{1\}} = \mathbf{\Phi}_{S, (\{1\} \cup S)^c} (\mathbf{\Phi}_{(\{1\} \cup S)^c, (\{1\} \cup S)^c})^{-1} \mathbf{\Phi}_{(\{1\} \cup S)^c, \{1\}}.$$

REMARK S1.1. Observe that for any $0 < \alpha < 1$, $g(\alpha, \gamma)$ is strictly increasing as a function of γ . Lemma S.1 therefore suggests that

$$g(\alpha, (\phi_0 \mathbf{e}_{1,p}^T \mathbf{\Phi}^{-1} \mathbf{e}_{1,p})^{-1} h^2) \leq g(\alpha, \{\phi_0 \mathbf{e}_{1,1+s}^T (\mathbf{\Phi}_{\{1\} \cup S, \{1\} \cup S})^{-1} \mathbf{e}_{1,1+s}\}^{-1} h^2),$$

with equality if and only if

$$\mathbf{\Phi}_{S, \{1\}} = \mathbf{\Phi}_{S, (\{1\} \cup S)^c} (\mathbf{\Phi}_{(\{1\} \cup S)^c, (\{1\} \cup S)^c})^{-1} \mathbf{\Phi}_{(\{1\} \cup S)^c, \{1\}}.$$

This implies that our test statistics are asymptotically more powerful than the Wald-type statistic based on the de-sparsified Lasso estimator.

S1.0.2. *Comparison with the decorrelated score test.* Let $\mathcal{M}_1 = [2, \dots, p]$. For any $\theta \in \mathbb{R}$, the decorrelated score function (Ning and Liu, 2017) is given by

$$\hat{S}_D(\theta) = \{\mathbf{X}_{\mathcal{M}_1} \hat{\boldsymbol{\omega}} - \mathbf{X}^1\}^T \{\mathbf{Y} - \mu(\theta \mathbf{X}^1 + \mathbf{X}_{\mathcal{M}_1} \hat{\boldsymbol{\beta}}_{\mathcal{M}_1})\},$$

for some penalized regression estimator $\hat{\boldsymbol{\beta}}$,

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \left(\frac{1}{n} \sum_{i=1}^n \{b(\boldsymbol{\beta}^T \mathbf{X}_i) - Y_i \boldsymbol{\beta}^T \mathbf{X}_i\} + \sum_{j=1}^p p_\lambda(|\beta_j|) \right),$$

and a Dantzig type estimator $\hat{\boldsymbol{\omega}}$,

$$\hat{\boldsymbol{\omega}} = \arg \min \|\boldsymbol{\omega}\|_1 \quad \text{s.t.} \quad \|(\boldsymbol{\omega}^T \mathbf{X}_{\mathcal{M}_1}^T - \mathbf{X}^1) \boldsymbol{\Sigma}(\mathbf{X} \hat{\boldsymbol{\beta}}) \mathbf{X}_{\mathcal{M}_1}\|_\infty \leq n\lambda'.$$

For testing $H_0 : \beta_{0,1} = 0$, [Ning and Liu \(2017\)](#) proposed the following decorrelated score test statistic,

$$T_{SD} = \frac{\sqrt{n} \hat{S}_D(0)}{\hat{\sigma}_s},$$

for some $\hat{\sigma}_s^2$ that consistently estimates $\sigma_s^2 = \phi_0(\mathbf{v}^*)^T \boldsymbol{\Phi} \mathbf{v}^*$ where $\mathbf{v}^* = (1, -(\boldsymbol{\omega}^*)^T)^T$ and $\boldsymbol{\omega}^* = (\boldsymbol{\Phi}_{\mathcal{M}_1, \mathcal{M}_1})^{-1} \boldsymbol{\Phi}_{\{1\}, \mathcal{M}_1}$. The null hypothesis is rejected when $|T_{SD}| > z_{\alpha/2}$. It was shown in Theorem 3.1 in [Ning and Liu \(2017\)](#) that such tests has correct size under H_0 .

Moreover, under $H_a : \beta_{0,1} = n^{-1/2}h$, it follows from Corollary D.1 in their supplementary article that, for any $0 < \alpha < 1$,

$$\begin{aligned} \Pr(|T_{SD}| > z_{\alpha/2}) &= \Pr\left(\mathbb{Z} > z_{\alpha/2} - h\sqrt{(\mathbf{v}^*)^T \boldsymbol{\Phi} \mathbf{v}^* / \phi_0}\right) \\ &- \Pr\left(\mathbb{Z} > z_{\alpha/2} + h\sqrt{(\mathbf{v}^*)^T \boldsymbol{\Phi} \mathbf{v}^* / \phi_0}\right) + o(1) = g(\alpha, h^2(\mathbf{v}^*)^T \boldsymbol{\Phi} \mathbf{v}^* / \phi_0) + o(1), \end{aligned}$$

where \mathbb{Z} denotes a standard normal random variable. By the definition of \mathbf{v}^* , we have

$$(\mathbf{v}^*)^T \boldsymbol{\Phi} \mathbf{v}^* = \boldsymbol{\Phi}_{\{1\}, \{1\}} - \boldsymbol{\Phi}_{\{1\}, \mathcal{M}_1} (\boldsymbol{\Phi}_{\mathcal{M}_1, \mathcal{M}_1})^{-1} \boldsymbol{\Phi}_{\mathcal{M}_1, \{1\}} = (\mathbf{e}_{1,p}^T \boldsymbol{\Phi}^{-1} \mathbf{e}_{1,p})^{-1},$$

where the last equality follows by the block matrix inversion formula (Lemma [S.9](#)). Therefore, it follows from Lemma [S.1](#) and the monotonicity of $g(\cdot)$ that our tests are more powerful than the decorrelated score test.

APPENDIX S2: ADDITIONAL DETAILS REGARDING $\hat{\phi}$ IN LINEAR REGRESSION MODELS

By Theorem [2.1](#), we have with probability tending to 1,

$$(S2.1) \quad \hat{S}_a = S.$$

Besides,

$$(S2.2) \quad \|\hat{\boldsymbol{\beta}}_{a, \text{SUM}} - \boldsymbol{\beta}_{a, \text{SUM}}\|_2 = O_p\left(\sqrt{(s+m)/n}\right).$$

Let $\varepsilon_i = Y_i - \mathbf{X}_i^T \boldsymbol{\beta}_0$. Under the event defined in (S2.1), we have

$$\begin{aligned} \hat{\phi} &= \underbrace{\frac{1}{n-s-m} \sum_{i=1}^n \varepsilon_i^2}_{I_1} + \underbrace{\frac{2}{n-s-m} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{i,\text{SUM}}^T (\hat{\boldsymbol{\beta}}_{a,\text{SUM}} - \boldsymbol{\beta}_{a,\text{SUM}})}_{I_2} \\ &+ \underbrace{\frac{1}{n-s-m} (\hat{\boldsymbol{\beta}}_{a,\text{SUM}} - \boldsymbol{\beta}_{0,\text{SUM}})^T \left(\sum_{i=1}^n \mathbf{X}_{i,\text{SUM}} \mathbf{X}_{i,\text{SUM}}^T \right) (\hat{\boldsymbol{\beta}}_{a,\text{SUM}} - \boldsymbol{\beta}_{0,\text{SUM}})}_{I_3}. \end{aligned}$$

By Chebyshev's inequality, we have for any $t_0 > 0$,

$$\Pr \left(\frac{1}{n} \sum_{i=1}^n (\varepsilon_i^2 - \phi_0) > t_0 \right) \leq \frac{\sum_{i=1}^n \text{Var}(\varepsilon_i^2)}{n^2 t_0^2}.$$

This implies

$$I_1 = \frac{n}{n-s-m} \phi_0 + O_p \left(\frac{\sqrt{n}}{n-s-m} \right).$$

Since we require $\max(s, m) = o(n)$, we have

$$(S2.3) \quad I_1 = \phi_0 + O_p(n^{-1/2}).$$

By Cauchy-Schwarz inequality,

$$|I_2| \leq \underbrace{\left\| \frac{2}{n-s-m} \sum_{i=1}^n \varepsilon_i \mathbf{X}_{i,\text{SUM}} \right\|_2}_{I_2^*} \|\hat{\boldsymbol{\beta}}_{a,\text{SUM}} - \boldsymbol{\beta}_{a,\text{SUM}}\|_2.$$

With some calculations, we have

$$\begin{aligned} \mathbb{E}(I_2^*)^2 &= \frac{4}{(n-s-m)^2} \mathbb{E} \sum_{i,j=1}^n \varepsilon_i \varepsilon_j \mathbf{X}_{i,\text{SUM}}^T \mathbf{X}_{j,\text{SUM}} \\ &= \frac{4\phi_0}{(n-s-m)^2} \sum_{i=1}^n \|\mathbf{X}_{i,\text{SUM}}\|_2^2 = \frac{\phi_0}{(n-s-m)^2} \text{tr}(\mathbf{X}_{\text{SUM}}^T \mathbf{X}_{\text{SUM}}) \\ &\leq \frac{4\phi_0(s+m)}{(n-s-m)^2} \lambda_{\max}(\mathbf{X}_{\text{SUM}}^T \mathbf{X}_{\text{SUM}}) = O \left(\frac{n(s+m)}{(n-s-m)^2} \right) = O \left(\frac{(s+m)}{n} \right), \end{aligned}$$

where the fourth equality is due to Condition (A1), and the last equality is due to the condition that $\max(s, m) = o(n)$. This implies

$$I_2^* = O_p \left(\frac{\sqrt{s+m}}{\sqrt{n}} \right),$$

which together with (S2.2) yields

$$(S2.4) \quad I_2 = O_p\left(\frac{s+m}{n}\right).$$

Besides, it follows (S2.2) and Condition (A1) that

$$\begin{aligned} I_3 &\leq \frac{1}{n-s-m} \|\hat{\beta}_{a,SUM} - \beta_{a,SUM}\|_2^2 \lambda_{\max}(\mathbf{X}_{SUM}^T \mathbf{X}_{SUM}) \\ &= O\left(\frac{s+m}{n-s-m}\right) = O\left(\frac{s+m}{n}\right). \end{aligned}$$

Combining this together with (S2.3) and (S2.4) gives that

$$\hat{\phi} = \phi_0 + O_p\left(\frac{1}{\sqrt{n}}\right) + O_p\left(\frac{s+m}{n}\right) = \phi_0 + O_p\left(\frac{1}{\sqrt{n}}\right),$$

where the last equality is due to the condition that $\max(s, m) = o(n^{1/2})$.

APPENDIX S3: ADDITIONAL DETAILS REGARDING POWERS OF PARTIAL PENALIZED TEST STATISTICS

The proposed test statistics are based on the partial penalized estimators $\hat{\beta}_0$, $\hat{\beta}_a$ defined in (2.2) and (2.3). Alternatively, we can construct our test statistics based on

$$\begin{aligned} \hat{\beta}_0^{\mathcal{N}} &= \arg \max_{\beta} Q_n^{\mathcal{N}}(\beta, \lambda_{n,0}) \text{ subject to } \mathbf{C}\beta_{\mathcal{M}} = \mathbf{t}, \\ \hat{\beta}_a^{\mathcal{N}} &= \arg \max_{\beta} Q_n^{\mathcal{N}}(\beta, \lambda_{n,a}), \end{aligned}$$

where

$$Q_n^{\mathcal{N}}(\beta, \lambda) = \frac{1}{n} \sum_{i=1}^n \{Y_i \beta^T \mathbf{X}_i - b(\beta^T \mathbf{X}_i)\} - \sum_{j \notin \mathcal{M} \cup \mathcal{N}} p\lambda(|\beta_j|),$$

for some set $\mathcal{N} \subseteq [1, \dots, p]$. The set \mathcal{N} can be chosen by the domain knowledge. For example, it can contain variables that the analyst thinks are important for ensuring robustness. Let $\mathcal{N}^* = \mathcal{N} \cap (\mathcal{M} \cup S)^c$. Define

$$\begin{aligned} T_L^{\mathcal{N}} &= 2n\{L_n(\hat{\beta}_a^{\mathcal{N}}) - L_n(\hat{\beta}_0^{\mathcal{N}})\}/\hat{\phi}, \\ T_W^{\mathcal{N}} &= (\mathbf{C}\hat{\beta}_{a,\mathcal{M}}^{\mathcal{N}} - \mathbf{t})^T (\mathbf{C}\hat{\Omega}_{a,mm}^{\mathcal{N}} \mathbf{C}^T)^{-1} (\mathbf{C}\hat{\beta}_{a,\mathcal{M}}^{\mathcal{N}} - \mathbf{t})/\hat{\phi}, \\ T_S^{\mathcal{N}} &= \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\hat{\beta}_0^{\mathcal{N}})\}^T \begin{pmatrix} \mathbf{X}_{\mathcal{M} \cup \mathcal{N}^*} \\ \mathbf{X}_{\hat{S}_0^{\mathcal{N}}} \end{pmatrix} \hat{\Omega}_0^{\mathcal{N}} \begin{pmatrix} \mathbf{X}_{\mathcal{M} \cup \mathcal{N}^*} \\ \mathbf{X}_{\hat{S}_0^{\mathcal{N}}} \end{pmatrix}^T \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\hat{\beta}_0^{\mathcal{N}})\}/\hat{\phi}, \end{aligned}$$

where $\hat{\phi}$ is some constant estimators for ϕ_0 , $\hat{\Omega}_{a,mm}^N$ is the first m rows and columns of $\hat{\Omega}_a^N$,

$$\begin{aligned}\hat{\Omega}_a^N &= n \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \Sigma(\mathbf{X} \hat{\beta}_a^N) \mathbf{X}_{\mathcal{M}} & \mathbf{X}_{\mathcal{M}}^T \Sigma(\mathbf{X} \hat{\beta}_a^N) \mathbf{X}_{\hat{\mathcal{S}}_a^N \cup \mathcal{N}^*} \\ \mathbf{X}_{\hat{\mathcal{S}}_a^N \cup \mathcal{N}^*}^T \Sigma(\mathbf{X} \hat{\beta}_a^N) \mathbf{X}_{\mathcal{M}} & \mathbf{X}_{\hat{\mathcal{S}}_a^N \cup \mathcal{N}^*}^T \Sigma(\mathbf{X} \hat{\beta}_a^N) \mathbf{X}_{\hat{\mathcal{S}}_a^N \cup \mathcal{N}^*} \end{pmatrix}^{-1}, \\ \hat{\Omega}_0^N &= n \begin{pmatrix} \mathbf{X}_{\mathcal{M} \cup \mathcal{N}^*}^T \Sigma(\mathbf{X} \hat{\beta}_0^N) \mathbf{X}_{\mathcal{M} \cup \mathcal{N}^*} & \mathbf{X}_{\mathcal{M} \cup \mathcal{N}^*}^T \Sigma(\mathbf{X} \hat{\beta}_0^N) \mathbf{X}_{\hat{\mathcal{S}}_0^N} \\ \mathbf{X}_{\hat{\mathcal{S}}_0^N}^T \Sigma(\mathbf{X} \hat{\beta}_0^N) \mathbf{X}_{\mathcal{M} \cup \mathcal{N}^*} & \mathbf{X}_{\hat{\mathcal{S}}_0^N}^T \Sigma(\mathbf{X} \hat{\beta}_0^N) \mathbf{X}_{\hat{\mathcal{S}}_0^N} \end{pmatrix}^{-1},\end{aligned}$$

and

$$\hat{\mathcal{S}}_a^N = \{j \in (\mathcal{M} \cup \mathcal{N}^*)^c : \hat{\beta}_{a,j}^N \neq 0\}, \quad \hat{\mathcal{S}}_0^N = \{j \in (\mathcal{M} \cup \mathcal{N}^*)^c : \hat{\beta}_{0,j}^N \neq 0\}.$$

For a given significance level α , we reject the null hypothesis when $T^N > \chi_\alpha^2(r)$ for $T^N = T_L^N, T_W^N$ or T_S^N .

Similar to Corollary 3.1, we can show the Type I error rates of T_L^N, T_W^N and T_S^N are close to the nominal level. Besides, under the alternative $\mathbf{C}\beta_{0,\mathcal{M}} - \mathbf{t} = \mathbf{h}_n$, the asymptotic power functions of these test statistics are equal to

$$(S3.1) \quad \Pr(\chi^2(r, \gamma_n^N) > \chi_\alpha^2(r)),$$

where $\gamma_n^N = n \mathbf{h}_n^T (\mathbf{C}\Omega_{mm}^N \mathbf{C}^T)^{-1} \mathbf{h}_n / \phi_0$, Ω_{mm}^N is the first m rows and columns of

$$\Omega_n^N = \left\{ \frac{1}{n} \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \Sigma(\mathbf{X} \beta_0) \mathbf{X}_{\mathcal{M}} & \mathbf{X}_{\mathcal{M}}^T \Sigma(\mathbf{X} \beta_0) \mathbf{X}_{\mathcal{S}} & \mathbf{X}_{\mathcal{M}}^T \Sigma(\mathbf{X} \beta_0) \mathbf{X}_{\mathcal{N}^*} \\ \mathbf{X}_{\mathcal{S}}^T \Sigma(\mathbf{X} \beta_0) \mathbf{X}_{\mathcal{M}} & \mathbf{X}_{\mathcal{S}}^T \Sigma(\mathbf{X} \beta_0) \mathbf{X}_{\mathcal{S}} & \mathbf{X}_{\mathcal{S}}^T \Sigma(\mathbf{X} \beta_0) \mathbf{X}_{\mathcal{N}^*} \\ \mathbf{X}_{\mathcal{N}^*}^T \Sigma(\mathbf{X} \beta_0) \mathbf{X}_{\mathcal{M}} & \mathbf{X}_{\mathcal{N}^*}^T \Sigma(\mathbf{X} \beta_0) \mathbf{X}_{\mathcal{S}} & \mathbf{X}_{\mathcal{N}^*}^T \Sigma(\mathbf{X} \beta_0) \mathbf{X}_{\mathcal{N}^*} \end{pmatrix} \right\}^{-1}.$$

We now show the test statistic achieves its greatest power if $\mathcal{N}^* = 0$. This means the partial penalized tests are most advantageous if each unpenalized variable is either an important variable (i.e., in \mathcal{S}) or a variable in \mathcal{M} . By (S3.1), it suffices to prove $\gamma_n^N \leq \gamma_n$ for any \mathcal{N} . Assume for now, we've shown

$$(S3.2) \quad \inf_{\mathbf{a} \in \mathbb{R}^m} (\mathbf{a}^T \Omega_{mm}^N \mathbf{a} - \mathbf{a}^T \Omega_{mm} \mathbf{a}) \geq 0.$$

This implies

$$(S3.3) \quad \inf_{\mathbf{a} \in \mathbb{R}^r} (\mathbf{a}^T \mathbf{C}\Omega_{mm}^N \mathbf{C}^T \mathbf{a} - \mathbf{a}^T \mathbf{C}\Omega_{mm} \mathbf{C}^T \mathbf{a}) \geq 0,$$

and hence the matrix $\mathbf{C}\Omega_{mm}^N \mathbf{C}^T - \mathbf{C}\Omega_{mm} \mathbf{C}^T$ is positive semidefinite.

Note that $\mathbf{C}\boldsymbol{\Omega}_{mm}\mathbf{C}^T$ is positive definite. By the eigenvalue decomposition theorem, we can find some positive definite matrix \mathbf{Q} such that $\mathbf{Q}\mathbf{Q} = \mathbf{C}\boldsymbol{\Omega}_{mm}\mathbf{C}^T$. By (S3.3), we can similarly show that the matrix

$$\mathbf{Q}^{-1}\mathbf{C}\boldsymbol{\Omega}_{mm}^{\mathcal{N}}\mathbf{C}^T\mathbf{Q}^{-1} - \mathbf{I}$$

is positive semidefinite. As a result, the eigenvalues of $\mathbf{Q}^{-1}\mathbf{C}\boldsymbol{\Omega}_{mm}^{\mathcal{N}}\mathbf{C}^T\mathbf{Q}^{-1}$ are all greater than or equal to 1. This implies that the eigenvalues of $\mathbf{Q}(\mathbf{C}\boldsymbol{\Omega}_{mm}^{\mathcal{N}}\mathbf{C}^T)^{-1}\mathbf{Q}^{-1}$ are all smaller than or equal to 1. Therefore, the matrix

$$\mathbf{I} - \mathbf{Q}(\mathbf{C}\boldsymbol{\Omega}_{mm}^{\mathcal{N}}\mathbf{C}^T)^{-1}\mathbf{Q}$$

is positive semidefinite. This further implies that the matrix $(\mathbf{C}\boldsymbol{\Omega}_{mm}\mathbf{C}^T)^{-1} - (\mathbf{C}\boldsymbol{\Omega}_{mm}^{\mathcal{N}}\mathbf{C}^T)^{-1}$ is positive semidefinite. Therefore, we have $\gamma_n \geq \gamma_n^{\mathcal{N}}$.

It remains to prove (S3.2). Let $\boldsymbol{\Omega}_{m+s, m+s}^{\mathcal{N}}$ be the first $(s+m)$ rows and columns of $\boldsymbol{\Omega}_n^{\mathcal{N}}$. By Lemma S.1, the matrix $\boldsymbol{\Omega}_{m+s, m+s}^{\mathcal{N}} - \boldsymbol{\Omega}_n$ is positive semidefinite. The assertion (S3.2) thus follows.

APPENDIX S4: DISCUSSION OF THE TECHNICAL CONDITIONS

S4.1. Discussion of Condition (A1). For Gaussian linear regression models, Condition (A1) reduces to (A1*) given below.

(A1*) Assume

$$\begin{aligned} \max_{1 \leq j \leq p} \|\mathbf{X}^j\|_{\infty} &= O\left(\sqrt{n/\log(p)}\right), & \max_{1 \leq j \leq p} \|\mathbf{X}^j\|_2 &= O(\sqrt{n}), \\ \lambda_{\min}(\mathbf{X}_{S \cup \mathcal{M}}^T \mathbf{X}_{S \cup \mathcal{M}}) &\geq cn, & \lambda_{\max}(\mathbf{X}_{S \cup \mathcal{M}}^T \mathbf{X}_{S \cup \mathcal{M}}) &= O(n), \\ \|\mathbf{X}_{(S \cup \mathcal{M})^c}^T \mathbf{X}_{S \cup \mathcal{M}}\|_{2, \infty} &= O(n). \end{aligned}$$

for some constants $c > 0$.

For logistic regression and Poisson regression models, Condition (A1) is implied by (A1*) and the following conditions:

$$(S4.1) \quad \max_{1 \leq i \leq n} |\mathbf{X}_i^T \boldsymbol{\beta}_0| = O(1),$$

$$(S4.2) \quad \max_{1 \leq j \leq p} \|\mathbf{X}^j\|_{\infty} = O\left(\frac{n^{1/2}}{(s+m)\log^{1/2} n}\right),$$

$$(S4.3) \quad \max_{1 \leq j \leq p} \lambda_{\max}\{\mathbf{X}_{S \cup \mathcal{M}}^T \text{diag}(|\mathbf{X}^j|) \mathbf{X}_{S \cup \mathcal{M}}\} = O(n).$$

We note that [van de Geer et al. \(2014\)](#) and [Ning and Liu \(2017\)](#) also assumed (S4.1) to establish the asymptotic properties of de-sparsified Lasso estimator

and the decorrelated score statistic under the high dimensional generalized linear models. Assume $\max(s, m) = O(n^{l_0})$ for some $0 < l_0 < 1/2$ and that each covariate has sub-exponential tail. Then we can show (S4.2) holds with probability tending to 1.

To further simplify Condition (A1*) and (S4.3), we consider a random design setting where $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ are i.i.d copies of some random vector $\mathbf{X}_0 \in \mathbb{R}^p$. More specifically, we consider the following two cases: (i) The bounded case: $\max_{j=[1, \dots, p]} |\mathbf{X}_{0,j}| \leq \omega_0$ for some constant $\omega_0 > 0$. (ii) The sub-Gaussian case: $\sup_{\|\mathbf{v}\|_2 \leq 1, \mathbf{v} \in \mathbb{R}^p} \|\mathbf{v}^T \mathbf{X}_0\|_{\psi_2} \leq \omega_0$ where for any random variable \mathbb{Z} ,

$$\|\mathbb{Z}\|_{\psi_2} \equiv \inf_{C>0} \left\{ \mathbb{E} \exp \left(\frac{|\mathbb{Z}|^2}{C^2} \right) \leq 2 \right\}.$$

Let $\mathbf{\Lambda} = \mathbb{E} \mathbf{X}_0 \mathbf{X}_0^T$.

LEMMA S.2. *Assume that $\lambda_{\max}(\mathbf{\Lambda}_{\text{SUM}, \text{SUM}}) = O(1)$, $\lambda_{\min}(\mathbf{\Lambda}_{\text{SUM}, \text{SUM}}) \geq \bar{c}$ for some constant $\bar{c} > 0$, $\max(s, m) = o\{(n/\log n)^{1/2}\}$ and $\max(s, m) \log p = o(n)$. Then in the bounded case, Condition (A1*) and Condition (S4.3) hold with probability tending to 1.*

LEMMA S.3. *Assume that $\lambda_{\min}(\mathbf{\Lambda}_{\text{SUM}, \text{SUM}}) \geq \bar{c}$ for some constant $\bar{c} > 0$, $\max(s, m) \log p = o(n/\log n)$ and $\log p = O(\sqrt{n})$. Then in the sub-Gaussian case, Condition (A1*) and Condition (S4.3) hold with probability tending to 1.*

Proofs of Lemma S.2 and Lemma S.3 are given in Section S5.3 and Section S5.4, respectively.

S4.2. Discussion of Condition (A2). The condition $p'_{\lambda_{n,j}}(d_n) = o((s+m)^{-1/2}n^{-1/2})$, $\lambda_{n,j}\kappa_{0,j} = o(1)$ automatically holds for SCAD penalty function, since $p'_{\lambda_{n,j}}(d_n) = 0$ and $\kappa_{0,j} = 0$ when $d_n \gg \lambda_{n,j}$. The Lasso penalty function doesn't satisfy the condition $p'_{\lambda_{n,j}}(d_n) = o((s+m)^{-1/2}n^{-1/2})$. Hence, the corresponding (un)constrained estimators have relatively large biases and don't have the asymptotic distributions in Theorem 2.1.

S4.3. Discussion of Condition (A3). For logistic regression models, we have $\max_i |Y_i - \mu(\mathbf{X}_i^T \boldsymbol{\beta}_0)| \leq 1$. It is easy to show Condition (A3) holds

with $M = 1$ and $v_0 = 2 \exp(1)$. For Gaussian linear models, we have

$$\begin{aligned} & \mathbb{E} \left\{ \exp \left(\frac{|Y_i - \mathbf{X}_i^T \boldsymbol{\beta}_0|}{\phi_0} \right) - 1 - \frac{|Y_i - \mathbf{X}_i^T \boldsymbol{\beta}_0|}{\phi_0} \right\} \\ & \leq \mathbb{E} \exp \left(\frac{|Y_i - \mathbf{X}_i^T \boldsymbol{\beta}_0|}{\phi_0} \right) \leq 2. \end{aligned}$$

Hence, (A3) holds with $M = \phi_0$ and $v_0 = 4\phi_0^2$. Assume

$$\max_{i \in [1, \dots, n]} |\mathbf{X}_i^T \boldsymbol{\beta}_0| \leq K_0,$$

for some constant $K_0 > 0$. Then, for Poisson regression models, we have

$$\begin{aligned} & \mathbb{E} \exp(|Y_i - \exp(\mathbf{X}_i^T \boldsymbol{\beta}_0)|) \leq \mathbb{E} \exp(|Y_i| + |\exp(\mathbf{X}_i^T \boldsymbol{\beta}_0)|) \\ & \leq \mathbb{E} \exp(2Y_i) = \exp \left\{ \exp(\mathbf{X}_i^T \boldsymbol{\beta}_0) (\exp(2) - 1) \right\} \leq \exp\{\exp(K_0 + 2)\}, \end{aligned}$$

where the second inequality follows by Jensen's inequality. Condition (A3) is thus satisfied.

S4.4. Discussion of Condition (A4). Condition (A4) is not restrictive at all. When the matrix \mathbf{C} doesn't vary with n , the condition $\lambda_{\max}((\mathbf{C}\mathbf{C}^T)^{-1}) = O(1)$ is automatically satisfied since \mathbf{C} is of full row rank.

APPENDIX S5: PROOFS

S5.1. Proof of Theorem 2.1. We focus on proving the statistical properties of $\hat{\boldsymbol{\beta}}_0$. Properties of $\hat{\boldsymbol{\beta}}_a$ can be similarly proven. We divide the proof into three steps. In the first step, we show that there exists a local maximizer $\hat{\boldsymbol{\beta}}$ of $Q_n(\boldsymbol{\beta})$ with the constraints $\mathbf{C}\boldsymbol{\beta}_{\mathcal{M}} = \mathbf{t}$ and $\boldsymbol{\beta}_{(\mathcal{M} \cup \mathcal{S})^c} = \mathbf{0}$, such that $\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0\|_2 = O_p(\sqrt{(s+m-r)/n})$. In the second step, we show $\hat{\boldsymbol{\beta}}$ is indeed a local maximizer of $Q_n(\boldsymbol{\beta})$ with the linear constraints $\mathbf{C}\boldsymbol{\beta}_{\mathcal{M}} = \mathbf{t}$. This implies $\hat{\boldsymbol{\beta}}_0 = \hat{\boldsymbol{\beta}}$. In the final step, we show

$$\begin{aligned} \sqrt{n} \begin{pmatrix} \hat{\boldsymbol{\beta}}_{0,\mathcal{M}} - \boldsymbol{\beta}_{0,\mathcal{M}} \\ \hat{\boldsymbol{\beta}}_{0,\mathcal{S}} - \boldsymbol{\beta}_{0,\mathcal{S}} \end{pmatrix} &= \frac{1}{\sqrt{n}} \mathbf{K}_n^{-1/2} (\mathbf{I} - \mathbf{P}_n) \mathbf{K}_n^{-1/2} \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_{\mathcal{S}}^T \end{pmatrix} \{ \mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}_0) \} \\ &\quad - \sqrt{n} \mathbf{K}_n^{-1/2} \mathbf{P}_n \mathbf{K}_n^{-1/2} \begin{pmatrix} \mathbf{C}^T (\mathbf{C}\mathbf{C}^T)^{-1} \mathbf{h}_n \\ \mathbf{0} \end{pmatrix} + o_p(1). \end{aligned}$$

Step 1: Define a p -dimensional vector $\boldsymbol{\beta}^*$ as

$$\begin{cases} \boldsymbol{\beta}_{\mathcal{M}}^* = \boldsymbol{\beta}_{0,\mathcal{M}} - \mathbf{C}^T (\mathbf{C}\mathbf{C}^T)^{-1} \mathbf{h}_n, \\ \boldsymbol{\beta}_{\mathcal{M}^c}^* = \boldsymbol{\beta}_{0,\mathcal{M}^c}. \end{cases}$$

It is immediate to see that

$$(S5.1) \quad \mathbf{C}\boldsymbol{\beta}_{\mathcal{M}}^* - \mathbf{t} = \mathbf{C}\boldsymbol{\beta}_{0,\mathcal{M}} - \mathbf{C}\mathbf{C}^T(\mathbf{C}\mathbf{C}^T)^{-1}\mathbf{h}_n - \mathbf{t} = \mathbf{h}_n - \mathbf{h}_n = \mathbf{0}.$$

Besides, it follows from Assumption (A4) that

$$(S5.2) \quad \begin{aligned} & \|\boldsymbol{\beta}_{\mathcal{M}\cup S}^* - \boldsymbol{\beta}_{0,\mathcal{M}\cup S}\|_2^2 \\ &= \|\mathbf{h}_n^T(\mathbf{C}\mathbf{C}^T)^{-1}\mathbf{h}_n\|_2^2 = O(\|\mathbf{h}_n\|_2^2) = O\left(\frac{m+s-r}{n}\right). \end{aligned}$$

Therefore, it suffices to show there exists a local maximizer $\hat{\boldsymbol{\beta}}$ of $Q_n(\boldsymbol{\beta})$ with the constraints $\mathbf{C}\boldsymbol{\beta}_{\mathcal{M}} = \mathbf{t}$ and $\boldsymbol{\beta}_{(\mathcal{M}\cup S)^c} = \mathbf{0}$, such that $\|\hat{\boldsymbol{\beta}}_{0,\mathcal{M}\cup S} - \boldsymbol{\beta}_{\mathcal{M}\cup S}^*\|_2^2 = O_p((m+s-r)/n)$.

Observe that for any $\boldsymbol{\beta}$ with $\mathbf{C}\boldsymbol{\beta}_{\mathcal{M}} = \mathbf{t}$, it follows from (S5.1) that $\mathbf{C}(\boldsymbol{\beta}_{\mathcal{M}} - \boldsymbol{\beta}_{\mathcal{M}}^*) = \mathbf{0}$ and hence $\boldsymbol{\beta}_{\mathcal{M}} - \boldsymbol{\beta}_{\mathcal{M}}^*$ belongs to the null space of \mathbf{C} . We take a basis matrix $\mathbf{Z} \in \mathbb{R}^{m \times (m-r)}$ of the null space \mathbf{C} . This implies $\mathbf{C}\mathbf{Z} = \mathbf{0}$. Further assume \mathbf{Z} is orthogonalized such that $\mathbf{Z}^T\mathbf{Z} = \mathbf{I}_{m-r}$ where \mathbf{I}_q stands for a $q \times q$ identity matrix. Hence, for any $\boldsymbol{\beta}_{\mathcal{M}}$ such that $\mathbf{C}\boldsymbol{\beta}_{\mathcal{M}} = \mathbf{t}$, there exists some $(m-r)$ -dimensional vector $\boldsymbol{\nu}$ such that $\boldsymbol{\beta}_{\mathcal{M}} - \boldsymbol{\beta}_{\mathcal{M}}^* = \mathbf{Z}\boldsymbol{\nu}$. For any $\boldsymbol{\delta} \in \mathbb{R}^{m-r+s}$, we define $\bar{Q}_n(\boldsymbol{\delta}) = Q_n(\boldsymbol{\beta}(\boldsymbol{\delta}))$ where $\boldsymbol{\beta}(\boldsymbol{\delta})$ is defined as

$$\begin{cases} \boldsymbol{\beta}(\boldsymbol{\delta})_{\mathcal{M}} = \boldsymbol{\beta}_{\mathcal{M}}^* + \mathbf{Z}\boldsymbol{\delta}_{J_0}, \\ \boldsymbol{\beta}(\boldsymbol{\delta})_S = \boldsymbol{\beta}_{0,S} + \boldsymbol{\delta}_{J_0^c}, \\ \boldsymbol{\beta}(\boldsymbol{\delta})_{(\mathcal{M}\cup S)^c} = \boldsymbol{\beta}_{0,(\mathcal{M}\cup S)^c}, \end{cases}$$

where $J_0 = [1, 2, \dots, m-r]$. Since $\|\mathbf{Z}\boldsymbol{\delta}_{J_0}\|_2^2 = \|\boldsymbol{\delta}_{J_0}\|_2^2$, it suffices to show that there exists a local maximizer $\boldsymbol{\delta}_0$ of $\bar{Q}_n(\boldsymbol{\delta})$ such that $\|\boldsymbol{\delta}_0\|_2 = O_p(\sqrt{(s+m-r)/n})$.

Define $\mathcal{N}_\tau = \{\boldsymbol{\delta} : \|\boldsymbol{\delta}\|_2 = \tau\}$, and

$$H_n = \left\{ \bar{Q}_n(\mathbf{0}) > \max_{\boldsymbol{\delta} \in \partial\mathcal{N}_\tau} \bar{Q}_n(\boldsymbol{\delta}) \right\},$$

where $\partial\mathcal{N}_\tau$ denotes the boundary of \mathcal{N}_τ . Clearly, on the event H_n , there exists a local maximizer $\boldsymbol{\delta}$ in \mathcal{N}_τ . Hence, it suffices to show $\Pr(H_n) \rightarrow 1$ as $n \rightarrow \infty$ and some sufficiently large τ .

For any $\boldsymbol{\delta}$, it follows from a second order Taylor's expansion that

$$(S5.3) \quad \bar{Q}_n(\boldsymbol{\delta}) - \bar{Q}_n(\mathbf{0}) = \boldsymbol{\delta}^T \mathbf{v} - \frac{1}{2} \boldsymbol{\delta}^T \mathbf{D} \boldsymbol{\delta},$$

where

$$\mathbf{v} = \frac{1}{n} \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix} \{ \mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}^*) \} - \begin{pmatrix} \mathbf{0}_{m-r} \\ \lambda_{n,0} \bar{\boldsymbol{\rho}}(\boldsymbol{\beta}_S^*) \end{pmatrix}$$

and

$$\mathbf{D} = \frac{1}{n} \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix} \Sigma(\mathbf{X}\boldsymbol{\beta}^{**}) \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix}^T - \begin{pmatrix} \mathbf{O}_{(m-r) \times (m-r)} & \mathbf{O}_{(m-r) \times s} \\ \mathbf{O}_{s \times (m-r)} & \boldsymbol{\Lambda}^* \end{pmatrix},$$

where $\boldsymbol{\beta}^{**}$ lies in the line segment jointing $\boldsymbol{\beta}(\boldsymbol{\delta})$ and $\boldsymbol{\beta}^*$, and $\boldsymbol{\Lambda}^*$ is a diagonal matrix with nonnegative diagonal elements. By the definition of $\boldsymbol{\beta}^*$ and $\boldsymbol{\beta}(\boldsymbol{\delta})$, we have $\boldsymbol{\beta}_{(\mathcal{M} \cup S)^c}^{**} = 0$. Moreover, it follows from (S5.2) and $\|\boldsymbol{\beta}(\boldsymbol{\delta}) - \boldsymbol{\beta}^*\|_2 = \tau \sqrt{(m+s-r)/n}$ that

$$\begin{aligned} \|\boldsymbol{\beta}^{**} - \boldsymbol{\beta}_0\|_2 &\leq \tau \sqrt{(s+m-r)/n} + O(\sqrt{(s+m-r)/n}) \\ &\ll \tau \sqrt{(s+m-r)/n} + \frac{1}{2} \sqrt{(s+m) \log(n)/n}. \end{aligned}$$

Therefore, for sufficiently large n and any $\tau \leq \sqrt{\log n}/2$, we have $\boldsymbol{\beta}^{**} \in \mathcal{N}_0$. By Condition (A2), the maximum eigenvalue of $\boldsymbol{\Lambda}^*$ is upper bounded by $\lambda_n \kappa_0$. Let

$$\mathbf{L} = \begin{pmatrix} \mathbf{Z} & \mathbf{O}_{m \times s} \\ \mathbf{O}_{s \times (m-r)} & \mathbf{I}_s \end{pmatrix},$$

for any $\boldsymbol{\delta} \in \mathbb{R}^{m+s-r}$, we have $\|\mathbf{L}\boldsymbol{\delta}\|_2 = \|\boldsymbol{\delta}\|_2$. Observe that

$$\frac{1}{n} \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix} \Sigma(\mathbf{X}\boldsymbol{\beta}^{**}) \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix}^T = \frac{1}{n} \mathbf{L}^T \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix} \Sigma(\mathbf{X}\boldsymbol{\beta}^{**}) \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix}^T \mathbf{L}.$$

Hence, for any $\boldsymbol{\delta} \in \mathbb{R}^{m+s-r}$ with $\|\boldsymbol{\delta}\|_2 = 1$, we have

$$\begin{aligned} &\frac{1}{n} \boldsymbol{\delta}^T \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix} \Sigma(\mathbf{X}\boldsymbol{\beta}^{**}) \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix}^T \boldsymbol{\delta} \\ &\geq \|\mathbf{L}\boldsymbol{\delta}\|_2^2 \lambda_{\min} \left\{ \frac{1}{n} \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix} \Sigma(\mathbf{X}\boldsymbol{\beta}^{**}) \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix}^T \right\} \\ &= \lambda_{\min} \left\{ \frac{1}{n} \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix} \Sigma(\mathbf{X}\boldsymbol{\beta}^{**}) \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix}^T \right\} \geq c, \end{aligned}$$

where the last inequality is due to Condition (A1). Since $\lambda_n \kappa_0 = o(1)$, we have $\lambda_{\min}(\mathbf{D}) \geq \bar{c}$ for some constant \bar{c} and sufficiently large n . Therefore, it follows from (S5.3) that

$$\begin{aligned} &\sup_{\boldsymbol{\delta} \in \mathcal{N}_\tau} \bar{Q}_n(\boldsymbol{\delta}) - \bar{Q}_n(\mathbf{0}) \leq \sup_{\boldsymbol{\delta} \in \mathcal{N}_\tau} \left(\|\boldsymbol{\delta}\|_2 \|\mathbf{v}\|_2 - \frac{\bar{c}}{2} \|\boldsymbol{\delta}\|_2^2 \right) \\ &= \tau \|\mathbf{v}\|_2 \sqrt{\frac{m+s-r}{n}} - \frac{\bar{c}\tau^2}{2} \frac{(m+s-r)}{n}. \end{aligned}$$

By Markov's inequality, for any fixed $\tau > 0$ and sufficiently large n , we have

$$(S5.4) \Pr(H_n) = \Pr\left(\|\mathbf{v}\|_2 < \bar{c}\tau\sqrt{\frac{m+s-r}{n}}\right) = 1 - \frac{n\mathbb{E}\|\mathbf{v}\|_2^2}{(m+s-r)\bar{c}^2\tau^2}.$$

Assume we can show

$$(S5.5) \quad \mathbb{E}\|\mathbf{v}\|_2^2 = O\left(\frac{m+s-r}{n}\right).$$

Let $\tau \rightarrow \infty$, it follows from (S5.4) that $\Pr(H_n) \rightarrow 1$ and the proof is hence completed. It remains to show (S5.5).

By the definition of \mathbf{v} , it follows from Cauchy-Schwarz inequality that

$$\begin{aligned} \mathbb{E}\|\mathbf{v}\|_2^2 &\leq 2\|\lambda_{n,0}\bar{\rho}(\boldsymbol{\beta}_S^*)\|_2^2 + 2\mathbb{E}\left\|\frac{1}{n}\begin{pmatrix} \mathbf{Z}^T \mathbf{X}_M^T \\ \mathbf{X}_S^T \end{pmatrix} \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}^*)\}\right\|_2^2 = 2\|\lambda_{n,0}\bar{\rho}(\boldsymbol{\beta}_S^*)\|_2^2 \\ &+ \frac{2}{n^2} \text{tr} \left[\mathbb{E} \left\{ \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_M^T \\ \mathbf{X}_S^T \end{pmatrix} \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}^*)\} \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}^*)\}^T \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_M^T \\ \mathbf{X}_S^T \end{pmatrix}^T \right\} \right] \\ &= 2\|\lambda_{n,0}\bar{\rho}(\boldsymbol{\beta}_S^*)\|_2^2 + \frac{2\phi_0}{n^2} \text{tr} \left\{ \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_M^T \\ \mathbf{X}_S^T \end{pmatrix} \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}_0) \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_M^T \\ \mathbf{X}_S^T \end{pmatrix}^T \right\} \\ &+ \frac{2}{n^2} \text{tr} \left\{ \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_M^T \\ \mathbf{X}_S^T \end{pmatrix} \{\boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}_0) - \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}^*)\} \{\boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}_0) - \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}^*)\}^T \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_M^T \\ \mathbf{X}_S^T \end{pmatrix}^T \right\} \\ &\triangleq 2I_0 + 2I_1 + 2I_2, \end{aligned}$$

where tr denotes the trace of a matrix. Observe that $\boldsymbol{\beta}_S^* = \boldsymbol{\beta}_{0,S}$. It follows from the monotonicity of ρ , the definition of d_n , and Condition (A2) that

$$I_0 \leq s\{\lambda_{n,0}\rho'(d_n)\}^2 = o\left(\frac{1}{n}\right).$$

Besides,

$$\begin{aligned} I_1 &\leq \frac{\phi_0(m+s-r)}{n^2} \lambda_{\max} \left\{ \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_M^T \\ \mathbf{X}_S^T \end{pmatrix} \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}_0) \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_M^T \\ \mathbf{X}_S^T \end{pmatrix}^T \right\} \\ &= \frac{\phi_0(m+s-r)}{n^2} \lambda_{\max} \left\{ \mathbf{L}^T \begin{pmatrix} \mathbf{X}_M^T \\ \mathbf{X}_S^T \end{pmatrix} \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}_0) \begin{pmatrix} \mathbf{X}_M^T \\ \mathbf{X}_S^T \end{pmatrix} \mathbf{L} \right\} \\ &= \frac{\phi_0(m+s-r)}{n^2} \lambda_{\max} \left\{ \begin{pmatrix} \mathbf{X}_M^T \\ \mathbf{X}_S^T \end{pmatrix} \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}_0) \begin{pmatrix} \mathbf{X}_M^T \\ \mathbf{X}_S^T \end{pmatrix}^T \right\} = O\left(\frac{s+m-r}{n}\right), \end{aligned}$$

where the last equality is due to Condition (A1).

To prove (S5.5), it remains to show $I_2 = O((m + s - r)/n)$. A first order Taylor expansion gives

$$\begin{aligned}\boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}_0) - \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}^*) &= \int_0^1 \boldsymbol{\Sigma}(\mathbf{X}\{\boldsymbol{\beta}^* + t(\boldsymbol{\beta}_0 - \boldsymbol{\beta}^*)\}) dt \mathbf{X}(\boldsymbol{\beta}_0 - \boldsymbol{\beta}^*) \\ &= \int_0^1 \boldsymbol{\Sigma}(\mathbf{X}\{\boldsymbol{\beta}^* + t(\boldsymbol{\beta}_0 - \boldsymbol{\beta}^*)\}) dt \mathbf{X}_{\mathcal{M}} \mathbf{C}^T (\mathbf{C} \mathbf{C}^T)^{-1} \mathbf{h}_n.\end{aligned}$$

Let $\mathcal{N}_0^* = \{\boldsymbol{\beta}^* + t(\boldsymbol{\beta}_0 - \boldsymbol{\beta}^*) : \forall 0 \leq t \leq 1\}$. By (S5.2), we have

$$\sup_{\boldsymbol{\beta} \in \mathcal{N}_0^*} \|\boldsymbol{\beta} - \boldsymbol{\beta}_0\|_2 = O\left(\frac{\sqrt{s + m - r}}{\sqrt{n}}\right).$$

It follows from Cauchy-Schwarz inequality that

$$\begin{aligned}& \left\| \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_{\mathcal{S}}^T \end{pmatrix} \int_0^1 \boldsymbol{\Sigma}(\mathbf{X}\{\boldsymbol{\beta}^* + t(\boldsymbol{\beta}_0 - \boldsymbol{\beta}^*)\}) dt \mathbf{X}_{\mathcal{M}} \right\|_2^2 \\ & \leq \sup_{\boldsymbol{\beta} \in \mathcal{N}_0^*} \lambda_{\max} \left\{ \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_{\mathcal{S}}^T \end{pmatrix} \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}) \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_{\mathcal{S}}^T \end{pmatrix}^T \right\} \\ & \times \sup_{\boldsymbol{\beta} \in \mathcal{N}_0^*} \lambda_{\max} (\mathbf{X}_{\mathcal{M}}^T \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}) \mathbf{X}_{\mathcal{M}}) \leq \sup_{\boldsymbol{\beta} \in \mathcal{N}_0^*} \left\{ \lambda_{\max} (\mathbf{X}_{\mathcal{MUS}}^T \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}) \mathbf{X}_{\mathcal{MUS}}) \right\}^2.\end{aligned}$$

By (A1), we have

$$(S5.6) \quad \lambda_{\max} (\mathbf{X}_{\mathcal{MUS}}^T \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}_0) \mathbf{X}_{\mathcal{MUS}}) = O(n).$$

Besides, it follows from Taylor's theorem that

$$\begin{aligned}& \sup_{\boldsymbol{\beta} \in \mathcal{N}_0^*} \sup_{\substack{\mathbf{a} \in \mathbb{R}^{s+m} \\ \|\mathbf{a}\|_2=1}} |\mathbf{a}^T \mathbf{X}_{\mathcal{MUS}}^T \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}) \mathbf{X}_{\mathcal{MUS}} \mathbf{a} - \mathbf{a}^T \mathbf{X}_{\mathcal{MUS}}^T \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}_0) \mathbf{X}_{\mathcal{MUS}} \mathbf{a}| \\ & \leq \sup_{\boldsymbol{\beta} \in \mathcal{N}_0^*} \sup_{\substack{\mathbf{a} \in \mathbb{R}^{s+m} \\ \|\mathbf{a}\|_2=1}} \sum_{j \in \mathcal{MUS}} |\mathbf{a}^T \mathbf{X}_{\mathcal{MUS}}^T \text{diag}\{|\mathbf{X}^j| \circ |b'''(\mathbf{X}\boldsymbol{\beta}^{**})|\} \mathbf{X}_{\mathcal{MUS}} \mathbf{a}| |\beta_{0,j} - \beta_j|,\end{aligned}$$

for some $\boldsymbol{\beta}^{**}$ lying on the line segment joining $\boldsymbol{\beta}_0$ and $\boldsymbol{\beta}$. For any $\boldsymbol{\beta} \in \mathcal{N}_0^*$, we have $\boldsymbol{\beta}^{**} \in \mathcal{N}_0^*$. Therefore,

$$\begin{aligned}& \sup_{\boldsymbol{\beta} \in \mathcal{N}_0^*} \sup_{\substack{\mathbf{a} \in \mathbb{R}^{s+m} \\ \|\mathbf{a}\|_2=1}} |\mathbf{a}^T \mathbf{X}_{\mathcal{MUS}}^T \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}) \mathbf{X}_{\mathcal{MUS}} \mathbf{a} - \mathbf{a}^T \mathbf{X}_{\mathcal{MUS}}^T \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}_0) \mathbf{X}_{\mathcal{MUS}} \mathbf{a}| \\ & \leq \sup_{\boldsymbol{\beta} \in \mathcal{N}_0^*} \sup_{\substack{\mathbf{a} \in \mathbb{R}^{s+m} \\ \|\mathbf{a}\|_2=1}} \sum_{j \in \mathcal{MUS}} |\mathbf{a}^T \mathbf{X}_{\mathcal{MUS}}^T \text{diag}\{|\mathbf{X}^j| \circ |b'''(\mathbf{X}\boldsymbol{\beta})|\} \mathbf{X}_{\mathcal{MUS}} \mathbf{a}| |\beta_{0,j} - \beta_j| \\ & \leq O(1)n \sup_{\boldsymbol{\beta} \in \mathcal{N}_0^*} \sum_j |\beta_{0,j} - \beta_j| = O(\sqrt{n}(s + m)) = O(n),\end{aligned}$$

where $O(1)$ denotes a positive constant, the second inequality is due to Condition (A1), the first equality follows by the definition of \mathcal{N}_0^* , and the last equality is due to the condition that $\max(s, m) = o(n^{1/2})$. Combining this together with (S5.6), we have

$$\begin{aligned}
\text{(S5.7)} \quad & \sup_{\beta \in \mathcal{N}_0^*} \lambda_{\max}(\mathbf{X}_{\mathcal{MUS}}^T \boldsymbol{\Sigma}(\mathbf{X}\beta) \mathbf{X}_{\mathcal{MUS}}) \\
&= \sup_{\beta \in \mathcal{N}_0^*} \sup_{\substack{\mathbf{a} \in \mathbb{R}^{s+m} \\ \|\mathbf{a}\|_2=1}} |\mathbf{a}^T \mathbf{X}_{\mathcal{MUS}}^T \boldsymbol{\Sigma}(\mathbf{X}\beta) \mathbf{X}_{\mathcal{MUS}} \mathbf{a}| \\
&\leq \sup_{\beta \in \mathcal{N}_0^*} \sup_{\substack{\mathbf{a} \in \mathbb{R}^{s+m} \\ \|\mathbf{a}\|_2=1}} |\mathbf{a}^T \mathbf{X}_{\mathcal{MUS}}^T \boldsymbol{\Sigma}(\mathbf{X}\beta) \mathbf{X}_{\mathcal{MUS}} \mathbf{a} - \mathbf{a}^T \mathbf{X}_{\mathcal{MUS}}^T \boldsymbol{\Sigma}(\mathbf{X}\beta_0) \mathbf{X}_{\mathcal{MUS}} \mathbf{a}| \\
&+ \sup_{\substack{\mathbf{a} \in \mathbb{R}^{s+m} \\ \|\mathbf{a}\|_2=1}} |\mathbf{a}^T \mathbf{X}_{\mathcal{MUS}}^T \boldsymbol{\Sigma}(\mathbf{X}\beta_0) \mathbf{X}_{\mathcal{MUS}} \mathbf{a}| = O(n),
\end{aligned}$$

and hence

$$\sup_{\beta \in \mathcal{N}_0^*} \{\lambda_{\max}(\mathbf{X}_{\mathcal{MUS}}^T \boldsymbol{\Sigma}(\mathbf{X}\beta) \mathbf{X}_{\mathcal{MUS}})\}^2 = O(n^2).$$

Using Cauchy-Schwarz inequality again, we obtain

$$\begin{aligned}
I_2 &\leq \frac{1}{n^2} \sup_{\beta \in \mathcal{N}_0^*} \left\| \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_{\mathcal{S}}^T \end{pmatrix} \boldsymbol{\Sigma}(\mathbf{X}\beta) \mathbf{X}_{\mathcal{M}} \right\|_2^2 \|\mathbf{C}^T (\mathbf{C}\mathbf{C}^T)^{-1} \mathbf{h}_n\|_2^2 \\
&= O(\|\mathbf{C}^T (\mathbf{C}\mathbf{C}^T)^{-1} \mathbf{h}_n\|_2^2) = O(\mathbf{h}_n^2) = O((s+m-r)/n),
\end{aligned}$$

by Conditions (A4). This completes the proof for the first step.

Step 2: In this step, we show that with probability tending to 1, the local maximizer $\hat{\beta}$ of $Q_n(\beta)$ with the constraints $\mathbf{C}\beta_{\mathcal{M}} = \mathbf{t}$ and $\beta_{(\mathcal{MUS})^c} = 0$ is indeed a local maximizer of $Q_n(\beta)$ with the linear constraints $\mathbf{C}\beta_{\mathcal{M}} = \mathbf{t}$. This implies $\hat{\beta}_0 = \hat{\beta}$. From the proof in the first step, we have shown that with probability at least $1 - 1/\tau^2$,

$$\text{(S5.8)} \quad \|\hat{\beta}_{\mathcal{MUS}} - \beta_{0,\mathcal{MUS}}\|_2 \leq \bar{c}\tau \left(\sqrt{\frac{s+m-r}{n}} \right), \quad \hat{\beta}_{(\mathcal{MUS})^c} = \mathbf{0},$$

for some constant $\bar{c} > 0$ and any $0 < \tau \leq \sqrt{\log n}/2$.

Similar to the proof of Theorem 1 in Fan and Lv (2011), it suffices to show with probability tending to 1, we have

$$\|\mathbf{X}_{(\mathcal{MUS})^c}^T \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\hat{\beta})\}\|_{\infty} < n\lambda_{n,0}\rho'(0+),$$

or

$$(S5.9) \quad \frac{1}{n} \|\mathbf{X}_{(\mathcal{M} \cup \mathcal{S})^c} \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\hat{\boldsymbol{\beta}})\}\|_\infty \ll \lambda_{n,0}.$$

Since $\hat{\boldsymbol{\beta}}_{(\mathcal{M} \cup \mathcal{S})^c} = \boldsymbol{\beta}_{0,(\mathcal{M} \cup \mathcal{S})^c} = \mathbf{0}$, using a second-order Taylor expansion around $\boldsymbol{\beta}_0$, we obtain for any $j \in (\mathcal{M} \cup \mathcal{S})^c$

$$(S5.10) \quad \begin{aligned} (\mathbf{X}^j)^T \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\hat{\boldsymbol{\beta}})\} &= (\mathbf{X}^j)^T \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}_0)\} \\ &- (\mathbf{X}^j)^T \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}_0) \mathbf{X}_{\mathcal{M} \cup \mathcal{S}} (\hat{\boldsymbol{\beta}}_{\mathcal{M} \cup \mathcal{S}} - \boldsymbol{\beta}_{0,\mathcal{M} \cup \mathcal{S}}) + R_j, \end{aligned}$$

where the absolute value of each element in the remainder term R_j is bounded by

$$\max_j (\hat{\boldsymbol{\beta}}_{\mathcal{M} \cup \mathcal{S}} - \boldsymbol{\beta}_{0,\mathcal{M} \cup \mathcal{S}})^T \mathbf{X}_{\mathcal{M} \cup \mathcal{S}}^T \text{diag}(|\mathbf{X}_j| \circ |b'''(\mathbf{X}\bar{\boldsymbol{\beta}}_j)|) \mathbf{X}_{\mathcal{M} \cup \mathcal{S}} (\hat{\boldsymbol{\beta}}_{\mathcal{M} \cup \mathcal{S}} - \boldsymbol{\beta}_{0,\mathcal{M} \cup \mathcal{S}}),$$

for some $\bar{\boldsymbol{\beta}}_j$ lying on the line segment joining $\boldsymbol{\beta}_0$ and $\hat{\boldsymbol{\beta}}$. Under the event defined in (S5.8), we have $\hat{\boldsymbol{\beta}} \in \mathcal{N}_0$. Since $\boldsymbol{\beta}_0 \in \mathcal{N}_0$, we have $\bar{\boldsymbol{\beta}}_j \in \mathcal{N}_0$. Therefore, under the event defined in (S5.8), it follows from Condition (A1) that

$$(S5.11) \quad \sup_{j \in (\mathcal{M} \cup \mathcal{S})^c} |R_j| \leq \bar{c}_1 \tau^2 (s+m) \ll \bar{c}_1 \tau^2 \sqrt{n(s+m)},$$

where \bar{c}_1 is some positive constant and the last equality is due to that $\max(s, m) = o(n^{1/2})$.

Besides, by Condition (A1) and (S5.8), there exists some constant $\bar{c}_2 > 0$ such that

$$(S5.12) \quad \begin{aligned} &\|\mathbf{X}_{(\mathcal{M} \cup \mathcal{S})^c}^T \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}_0) \mathbf{X}_{\mathcal{M} \cup \mathcal{S}} (\hat{\boldsymbol{\beta}}_{\mathcal{M} \cup \mathcal{S}} - \boldsymbol{\beta}_{0,\mathcal{M} \cup \mathcal{S}})\|_\infty \\ &\leq \|\mathbf{X}_{(\mathcal{M} \cup \mathcal{S})^c}^T \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}_0) \mathbf{X}_{\mathcal{M} \cup \mathcal{S}}\|_{2,\infty} \|\hat{\boldsymbol{\beta}}_{\mathcal{M} \cup \mathcal{S}} - \boldsymbol{\beta}_{0,\mathcal{M} \cup \mathcal{S}}\|_2 \leq \bar{c}_2 \tau \sqrt{n(s+m)}, \end{aligned}$$

with probability at least $1 - 1/\tau^2$.

Moreover, it follows from the condition $\max_{1 \leq j \leq p} \|\mathbf{X}^j\|_\infty = O(\sqrt{n/\log(p)})$, $\max_{1 \leq j \leq p} \|\mathbf{X}^j\|_2 = O(\sqrt{n})$ in (A1), Condition (A3) and Proposition 4 in Fan and Lv (2011) that

$$\begin{aligned} &\max_j \Pr \left(|(\mathbf{X}^j)^T \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}_0)\}| > \gamma \sqrt{n \log p} \right) \\ &\leq 2 \exp \left(-\frac{1}{2} \frac{\gamma^2 n \log p}{O(nv_0) + O(\sqrt{n/\log p}) M \gamma \sqrt{n \log p}} \right) \\ &\leq 2 \exp \left(-\frac{1}{2} \frac{\gamma^2 \log p}{O(v_0) + O(\gamma M)} \right), \end{aligned}$$

where the constants M and v_0 are defined in Condition (A3). By Bonferroni's inequality, we have

(S5.13)

$$\Pr \left(\max_j |X_j^T \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}_0)\}| > \gamma \sqrt{n \log p} \right) \leq 2 \exp \left(-\frac{1}{2} \frac{\gamma^2 \log p}{O(v_0 + \gamma M)} + \log p \right).$$

Since $p \rightarrow \infty$, for sufficiently large γ , the RHS of (S5.13) converges to 0. This implies we have with probability tending to 1,

$$(S5.14) \quad \|\mathbf{X}_{(\mathcal{MUS})^c}^T \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}_0)\}\|_\infty \leq \bar{c}_3 \sqrt{n \log p},$$

for some constant $\bar{c}_3 > 0$.

Let $\bar{c}_4 = \max(\bar{c}_1, \bar{c}_2, \bar{c}_3)$. Combining (S5.11), (S5.12) with (S5.14), we obtain

$$(S5.15) \quad \begin{aligned} \Pr \left(\|\mathbf{X}_{(\mathcal{MUS})^c}^T \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\hat{\boldsymbol{\beta}})\}\|_\infty \leq \bar{c}_4 \sqrt{n \log p} + \bar{c}_4 (\tau + \tau^2) \sqrt{n(s+m)} \right) \\ \geq 1 - \frac{1}{\tau^2} + o(1). \end{aligned}$$

By (A2), we have $\lambda_{n,0} \gg \max(\sqrt{\log p}, \sqrt{s+m})/\sqrt{n}$. Let τ_n to be any diverging sequence that $\tau_n \ll \sqrt{\log n}$ and $\tau_n \ll \sqrt{n}\lambda_{n,0}/\max(\sqrt{\log p}, \sqrt{s+m})$. By (S5.15), we have

$$\Pr \left(\|\mathbf{X}_{(\mathcal{MUS})^c}^T \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\hat{\boldsymbol{\beta}})\}\|_\infty \ll \lambda_{n,0} \right) \geq 1 - \frac{1}{\tau_n^2} + o(1) \rightarrow 1.$$

Thus, we have with probability tending to 1, (S5.9) holds.

Step 3: We've shown that for any $0 < \tau \leq \sqrt{\log n}/2$ and sufficiently large n ,

$$(S5.16) \quad \Pr \left\{ \|\hat{\boldsymbol{\beta}}_{0,\mathcal{MUS}} - \boldsymbol{\beta}_{0,\mathcal{MUS}}\|_2 \leq \bar{c}\tau \left(\frac{s+m-r}{n} \right) \right\} \geq 1 - \frac{1}{\tau^2},$$

$$(S5.17) \quad \Pr \left(\hat{\boldsymbol{\beta}}_0 \in \mathcal{N}_0 \right) \rightarrow 1,$$

$$(S5.18) \quad \mathbf{C}\hat{\boldsymbol{\beta}}_{0,\mathcal{M}} = \mathbf{t}.$$

Finally, we show that (S5.19) holds.

$$(S5.19) \quad \begin{aligned} \sqrt{n} \begin{pmatrix} \hat{\boldsymbol{\beta}}_{0,\mathcal{M}} - \boldsymbol{\beta}_{0,\mathcal{M}} \\ \hat{\boldsymbol{\beta}}_{0,\mathcal{S}} - \boldsymbol{\beta}_{0,\mathcal{S}} \end{pmatrix} &= \frac{1}{\sqrt{n}} \mathbf{K}_n^{-1/2} (\mathbf{I} - \mathbf{P}_n) \mathbf{K}_n^{-1/2} \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_{\mathcal{S}}^T \end{pmatrix} \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}_0)\} \\ &- \sqrt{n} \mathbf{K}_n^{-1/2} \mathbf{P}_n \mathbf{K}_n^{-1/2} \begin{pmatrix} \mathbf{C}^T (\mathbf{C}\mathbf{C}^T)^{-1} \mathbf{h}_n \\ \mathbf{0} \end{pmatrix} + o_p(1). \end{aligned}$$

Theorem 2.1 therefore follows.

In the first step, we have shown that $\hat{\beta}_0$ is the local maximizer of $Q_n(\beta)$ with the constraints $\mathbf{C}\beta_{\mathcal{M}} = \mathbf{t}$, $\beta_{(\mathcal{M} \cup \mathcal{S})^c} = \mathbf{0}$. This implies that there exists some vector $\nu \in \mathbb{R}^r$ such that

$$(S5.20) \quad \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_{\mathcal{S}}^T \end{pmatrix} \{ \mathbf{Y} - \mu(\mathbf{X}_{\mathcal{M} \cup \mathcal{S}} \hat{\beta}_{0, \mathcal{M} \cup \mathcal{S}}) \} = \begin{pmatrix} \sqrt{n} \mathbf{C}^T \nu \\ n \lambda_n \bar{\rho}(\hat{\beta}_{0, \mathcal{S}}) \end{pmatrix}.$$

Similar to (S5.10) and (S5.11), we can show that the left-hand side (LHS) of (S5.20) is equal to

$$(S5.21) \quad \begin{aligned} & \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_{\mathcal{S}}^T \end{pmatrix} \{ \mathbf{Y} - \mu(\mathbf{X}_{\mathcal{M} \cup \mathcal{S}} \beta_{0, \mathcal{M} \cup \mathcal{S}}) \} \\ & - \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_{\mathcal{S}}^T \end{pmatrix} \Sigma(\mathbf{X} \beta_0) \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_{\mathcal{S}}^T \end{pmatrix}^T \begin{pmatrix} \hat{\beta}_{0, \mathcal{M}} - \beta_{0, \mathcal{M}} \\ \hat{\beta}_{0, \mathcal{S}} - \beta_{0, \mathcal{S}} \end{pmatrix} + \mathbf{R}, \end{aligned}$$

where the remainder term \mathbf{R} satisfies $\|\mathbf{R}\|_{\infty} = O_p(s + m)$. Hence we have

$$\|\mathbf{R}\|_2 = O_p\left((s + m)^{3/2}\right) = o_p(\sqrt{n}),$$

under the conditions that $s + m = o(n^{1/3})$. Recall that

$$\mathbf{K}_n = \frac{1}{n} \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_{\mathcal{S}}^T \end{pmatrix} \Sigma(\mathbf{X} \beta_0) \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_{\mathcal{S}}^T \end{pmatrix}^T.$$

Under Condition (A1), we have $\liminf_n \lambda_{\min}(\mathbf{K}_n) > 0$. This implies $\|\mathbf{K}_n^{-1} \mathbf{R}\|_2 = o_p(\sqrt{n})$. Combining this together with (S5.20) and (S5.21), we obtain

$$(S5.22) \quad \begin{aligned} \sqrt{n} \begin{pmatrix} \hat{\beta}_{0, \mathcal{M}} - \beta_{0, \mathcal{M}} \\ \hat{\beta}_{0, \mathcal{S}} - \beta_{0, \mathcal{S}} \end{pmatrix} &= \frac{1}{\sqrt{n}} \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_{\mathcal{S}}^T \end{pmatrix} \{ \mathbf{Y} - \mu(\mathbf{X}_{\mathcal{M} \cup \mathcal{S}} \beta_{0, \mathcal{M} \cup \mathcal{S}}) \} \\ &- \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{C}^T \nu \\ \sqrt{n} \lambda_n \bar{\rho}(\hat{\beta}_{0, \mathcal{S}}) \end{pmatrix} + o_p(1). \end{aligned}$$

Besides, by (S5.16), it follows from Condition (A2) that $\|\hat{\beta}_{0, \mathcal{S}} - \beta_{0, \mathcal{S}}\|_{\infty} \ll d_n$ with probability tending to 1. This implies for sufficiently large n , we have $\min_{j \in \mathcal{S}} |\hat{\beta}_{0, j}| > \min_{j \in \mathcal{S}} |\beta_{0, j}| - d_n = d_n$. By the monotonicity of ρ' and Condition (A2), we obtain

$$|\sqrt{n} \lambda_n \bar{\rho}(\hat{\beta}_{0, \mathcal{S}})| \leq s^{1/2} |\sqrt{n} \lambda_n \rho'(d_n)| = o(1),$$

with probability tending to 1. This together with (S5.22) suggests that

$$(S5.23) \quad \begin{aligned} \sqrt{n} \begin{pmatrix} \hat{\beta}_{0,\mathcal{M}} - \beta_{0,\mathcal{M}} \\ \hat{\beta}_{0,S} - \beta_{0,S} \end{pmatrix} &= \frac{1}{\sqrt{n}} \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix} \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}_{\mathcal{M} \cup S} \boldsymbol{\beta}_{0,\mathcal{M} \cup S})\} \\ &- \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix} \boldsymbol{\nu} + o_p(1). \end{aligned}$$

Since $\mathbf{C}\boldsymbol{\beta}_{0,\mathcal{M}} - \mathbf{t} = \mathbf{h}_n$, by (S5.18), we have $\mathbf{C}(\hat{\beta}_{0,\mathcal{M}} - \beta_{0,\mathcal{M}}) = -\mathbf{h}_n$. Therefore, it follows from (S5.23) that

$$(S5.24) \quad \begin{aligned} &\begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix}^T \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix} \boldsymbol{\nu} - \sqrt{n} \mathbf{h}_n \\ &= \frac{1}{\sqrt{n}} \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix}^T \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix} \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}_{\mathcal{M} \cup S} \boldsymbol{\beta}_{0,\mathcal{M} \cup S})\} + \mathbf{C}\mathbf{R}^*, \end{aligned}$$

for some m -dimensional vector \mathbf{R}^* such that $\|\mathbf{R}^*\|_2 = o_p(1)$.

By definition, we have

$$\begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix}^T \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix} = \mathbf{C}^T \boldsymbol{\Omega}_{mm} \mathbf{C}.$$

Multiplying $(\mathbf{C}\boldsymbol{\Omega}_{mm}\mathbf{C}^T)^{-1}$ on both sides of (S5.24), we have

$$\begin{aligned} \boldsymbol{\nu} &= \frac{1}{\sqrt{n}} \boldsymbol{\Psi}^{-1} \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix}^T \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix} \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}_{\mathcal{M} \cup S} \boldsymbol{\beta}_{0,\mathcal{M} \cup S})\} \\ &+ \boldsymbol{\Psi}^{-1} \mathbf{C}\mathbf{R}^* + \sqrt{n} \boldsymbol{\Psi}^{-1} \mathbf{h}_n, \end{aligned}$$

where $\boldsymbol{\Psi} = \mathbf{C}\boldsymbol{\Omega}_{mm}\mathbf{C}^T$. This together with (S5.23) yields

$$(S5.25) \quad \begin{aligned} &\sqrt{n} \begin{pmatrix} \hat{\beta}_{0,\mathcal{M}} - \beta_{0,\mathcal{M}} \\ \hat{\beta}_{0,S} - \beta_{0,S} \end{pmatrix} = \frac{1}{\sqrt{n}} \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix} \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}_{\mathcal{M} \cup S} \boldsymbol{\beta}_{0,\mathcal{M} \cup S})\} \\ &- \frac{1}{\sqrt{n}} \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix} \boldsymbol{\Psi}^{-1} \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix}^T \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix} \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}_{\mathcal{M} \cup S} \boldsymbol{\beta}_{0,\mathcal{M} \cup S})\} \\ &- \sqrt{n} \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix} \boldsymbol{\Psi}^{-1} \mathbf{h}_n - \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix} \boldsymbol{\Psi}^{-1} \mathbf{C}\mathbf{R}^* + o_p(1) \\ &= \frac{1}{\sqrt{n}} \mathbf{K}_n^{-1/2} (\mathbf{I} - \mathbf{P}_n) \mathbf{K}_n^{-1/2} \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix} \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}_{\mathcal{M} \cup S} \boldsymbol{\beta}_{0,\mathcal{M} \cup S})\} \\ &- \sqrt{n} \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix} \boldsymbol{\Psi}^{-1} \mathbf{h}_n - \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix} \boldsymbol{\Psi}^{-1} \mathbf{C}\mathbf{R}^* + o_p(1). \end{aligned}$$

In the following, we prove

$$(S5.26) \quad \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix} \Psi^{-1} \mathbf{C} \mathbf{R}^* = \boldsymbol{\Omega}_{mm} \mathbf{C}^T \Psi^{-1} \mathbf{C} \mathbf{R}^* = o_p(1).$$

Under Condition (A1), we have $\lambda_{\max}(\mathbf{K}_n) = O(1)$. This implies $\liminf_n \lambda_{\min}(\boldsymbol{\Omega}_n) > 0$, or equivalently,

$$\inf_{\mathbf{a} \in \mathbb{R}^{m+s}: \|\mathbf{a}\|_2=1} \liminf_n \mathbf{a}^T \boldsymbol{\Omega}_n \mathbf{a} > 0.$$

Hence, we have

$$\inf_{\mathbf{a} \in \mathbb{R}^{m+s}: \|\mathbf{a}\|_2=1, \mathbf{a}_{J_0}^c=0} \liminf_n \mathbf{a}^T \boldsymbol{\Omega}_n \mathbf{a} > 0,$$

where $J_0 = [1, \dots, m]$. Note that this implies

$$\inf_{\mathbf{a} \in \mathbb{R}^{m+s}: \|\mathbf{a}\|_2=1} \liminf_n \mathbf{a}^T \boldsymbol{\Omega}_{mm} \mathbf{a} > 0.$$

Therefore, we obtain

$$(S5.27) \quad \liminf_n \lambda_{\min}(\boldsymbol{\Omega}_{mm}) > 0.$$

Similarly, we can show

$$(S5.28) \quad \lambda_{\max}(\boldsymbol{\Omega}_{mm}) = O(1).$$

Using Cauchy-Schwarz inequality, we have

$$\begin{aligned} & \|\boldsymbol{\Omega}_{mm} \mathbf{C}^T (\mathbf{C} \boldsymbol{\Omega}_{mm} \mathbf{C}^T)^{-1} \mathbf{C} \mathbf{R}^*\|_2^2 \\ &= \|\boldsymbol{\Omega}_{mm}^{1/2} \boldsymbol{\Omega}_{mm}^{1/2} \mathbf{C}^T (\mathbf{C} \boldsymbol{\Omega}_{mm} \mathbf{C}^T)^{-1} \mathbf{C} \boldsymbol{\Omega}_{mm}^{1/2} \boldsymbol{\Omega}_{mm}^{-1/2} \mathbf{R}^*\|_2^2 \\ &\leq \|\boldsymbol{\Omega}_{mm}^{1/2}\|_2^2 \|\boldsymbol{\Omega}_{mm}^{1/2} \mathbf{C}^T (\mathbf{C} \boldsymbol{\Omega}_{mm} \mathbf{C}^T)^{-1} \mathbf{C} \boldsymbol{\Omega}_{mm}^{1/2}\|_2^2 \|\boldsymbol{\Omega}_{mm}^{-1/2}\|_2^2 \|\mathbf{R}^*\|_2^2 \\ &\leq \lambda_{\max}(\boldsymbol{\Omega}_{mm}) \lambda_{\max}(\boldsymbol{\Omega}_{mm}^{-1}) \|\mathbf{R}^*\|_2^2 = o_p(1), \end{aligned}$$

by (S5.27) and (S5.28). This proves (S5.26). Hence, it follows from (S5.25) that

$$(S5.29) \quad \begin{aligned} \sqrt{n} \begin{pmatrix} \hat{\boldsymbol{\beta}}_{0,\mathcal{M}} - \boldsymbol{\beta}_{0,\mathcal{M}} \\ \hat{\boldsymbol{\beta}}_{0,S} - \boldsymbol{\beta}_{0,S} \end{pmatrix} &= \frac{1}{\sqrt{n}} \mathbf{K}_n^{-1/2} (\mathbf{I} - \mathbf{P}_n) \mathbf{K}_n^{-1/2} \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix} \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X} \boldsymbol{\beta}_0)\} \\ &- \sqrt{n} \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix} \Psi^{-1} \mathbf{h}_n + o_p(1). \end{aligned}$$

Moreover, observe that

$$\mathbf{h}_n = \mathbf{C}\mathbf{C}^T(\mathbf{C}\mathbf{C}^T)^{-1}\mathbf{h}_n = \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix}^T \begin{pmatrix} \mathbf{C}^T(\mathbf{C}\mathbf{C}^T)^{-1}\mathbf{h}_n \\ \mathbf{0} \end{pmatrix}.$$

Hence, we have

$$\begin{aligned} \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix} \boldsymbol{\Psi}^{-1}\mathbf{h}_n &= \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix} \boldsymbol{\Psi}^{-1} \begin{pmatrix} \mathbf{C}^T \\ \mathbf{O}_{r \times s}^T \end{pmatrix}^T \begin{pmatrix} \mathbf{C}^T(\mathbf{C}\mathbf{C}^T)^{-1}\mathbf{h}_n \\ \mathbf{0} \end{pmatrix} \\ &= \mathbf{K}_n^{-1/2} \mathbf{P}_n \mathbf{K}_n^{-1/2} \begin{pmatrix} \mathbf{C}^T(\mathbf{C}\mathbf{C}^T)^{-1}\mathbf{h}_n \\ \mathbf{0} \end{pmatrix}. \end{aligned}$$

This together with (S5.29) proves (S5.19).

Similarly, we can show $\hat{\boldsymbol{\beta}}_a$ satisfies (i) $\Pr(\hat{\boldsymbol{\beta}}_{a,(\mathcal{M} \cup S)^c} = \mathbf{0}) \rightarrow 1$; (ii) $\|\hat{\boldsymbol{\beta}}_{a,\mathcal{M} \cup S} - \boldsymbol{\beta}_{0,\mathcal{M} \cup S}\|_2 = O_p(\sqrt{(s+m)/n})$; (iii)

$$\begin{aligned} \sqrt{n} \begin{pmatrix} \hat{\boldsymbol{\beta}}_{a,\mathcal{M}} - \boldsymbol{\beta}_{0,\mathcal{M}} \\ \hat{\boldsymbol{\beta}}_{a,S} - \boldsymbol{\beta}_{0,S} \end{pmatrix} &= \frac{1}{\sqrt{n}} \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix} \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}_0)\} + O\left(\frac{(s+m)^{3/2}}{\sqrt{n}}\right) \\ &+ O_p\left((ns)^{1/2}\lambda_{n,a}\rho'(d_n, \lambda_{n,a})\right). \end{aligned}$$

Under the condition that $\max(s, m) = o(n^{1/3})$ and $\lambda_{n,a}\rho'(d_n, \lambda_{n,a}) = o(n^{-1/2}(s+m)^{-1/2})$ in (A2), we obtain

$$\sqrt{n} \begin{pmatrix} \hat{\boldsymbol{\beta}}_{a,\mathcal{M}} - \boldsymbol{\beta}_{0,\mathcal{M}} \\ \hat{\boldsymbol{\beta}}_{a,S} - \boldsymbol{\beta}_{0,S} \end{pmatrix} = \frac{1}{\sqrt{n}} \mathbf{K}_n^{-1} \begin{pmatrix} \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_S^T \end{pmatrix} \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}_0)\} + o_p(1).$$

The proof is hence completed.

S5.2. Proof of Lemma S.1. Let $\boldsymbol{\Theta} = \boldsymbol{\Phi}^{-1}$. It follows from the block matrix inversion formula (Lemma S.9) that

$$\begin{aligned} &\boldsymbol{\Theta}_{\{1\} \cup S, \{1\} \cup S}^{-1} \\ &= \boldsymbol{\Phi}_{\{1\} \cup S, \{1\} \cup S} - \boldsymbol{\Phi}_{\{1\} \cup S, (\{1\} \cup S)^c} (\boldsymbol{\Phi}_{(\{1\} \cup S)^c, (\{1\} \cup S)^c})^{-1} \boldsymbol{\Phi}_{(\{1\} \cup S)^c, \{1\} \cup S}. \end{aligned}$$

Observe that $\mathbf{e}_{1,p}^T \boldsymbol{\Phi}^{-1} \mathbf{e}_{1,p} = \mathbf{e}_{1,p}^T \boldsymbol{\Theta} \mathbf{e}_{1,p} = \mathbf{e}_{1,1+s}^T \boldsymbol{\Theta}_{1 \cup \{S\}, 1 \cup \{S\}} \mathbf{e}_{1,1+s}$. Hence,

$$\begin{aligned} \mathbf{e}_{1,p}^T \boldsymbol{\Phi}^{-1} \mathbf{e}_{1,p} &= \mathbf{e}_{1,1+s}^T \\ &\left\{ \boldsymbol{\Phi}_{\{1\} \cup S, \{1\} \cup S} - \boldsymbol{\Phi}_{\{1\} \cup S, (\{1\} \cup S)^c} (\boldsymbol{\Phi}_{(\{1\} \cup S)^c, (\{1\} \cup S)^c})^{-1} \boldsymbol{\Phi}_{(\{1\} \cup S)^c, \{1\} \cup S} \right\}^{-1} \mathbf{e}_{1,1+s}. \end{aligned}$$

With some calculation, we have

$$\begin{aligned}
& \Phi_{\{1\} \cup S, \{1\} \cup S} - \Phi_{\{1\} \cup S, (\{1\} \cup S)^c} (\Phi_{(\{1\} \cup S)^c, (\{1\} \cup S)^c})^{-1} \Phi_{(\{1\} \cup S)^c, \{1\} \cup S} \\
= & \begin{pmatrix} \Phi_{\{1\}, \{1\}} & \Phi_{S, \{1\}}^T \\ \Phi_{S, \{1\}} & \Phi_{S, S} \end{pmatrix} - \begin{pmatrix} \Phi_{\{1\}, (\{1\} \cup S)^c} \\ \Phi_{S, (\{1\} \cup S)^c} \end{pmatrix} \Phi_{(\{1\} \cup S)^c, (\{1\} \cup S)^c}^{-1} \begin{pmatrix} \Phi_{\{1\}, (\{1\} \cup S)^c} \\ \Phi_{S, (\{1\} \cup S)^c} \end{pmatrix}^T \\
= & \begin{pmatrix} \mathbf{K}_{11} & \mathbf{K}_{21}^T \\ \mathbf{K}_{21} & \mathbf{K}_{22} \end{pmatrix},
\end{aligned}$$

where

$$\begin{aligned}
\mathbf{K}_{11} &= \Phi_{\{1\}, \{1\}} - \Phi_{\{1\}, (\{1\} \cup S)^c} (\Phi_{(\{1\} \cup S)^c, (\{1\} \cup S)^c})^{-1} \Phi_{\{1\}, (\{1\} \cup S)^c}^T, \\
\mathbf{K}_{21} &= \Phi_{S, \{1\}} - \Phi_{S, (\{1\} \cup S)^c} (\Phi_{(\{1\} \cup S)^c, (\{1\} \cup S)^c})^{-1} \Phi_{(\{1\} \cup S)^c, \{1\}}, \\
\mathbf{K}_{22} &= \Phi_{S, S} - \Phi_{S, (\{1\} \cup S)^c} (\Phi_{(\{1\} \cup S)^c, (\{1\} \cup S)^c})^{-1} \Phi_{(\{1\} \cup S)^c, S}.
\end{aligned}$$

By Lemma S.9, we have

$$e_{1,p}^T \Phi^{-1} e_{1,p} = (\mathbf{K}_{11} - \mathbf{K}_{21}^T \mathbf{K}_{22}^{-1} \mathbf{K}_{21})^{-1}.$$

Similarly we have

$$e_{1,1+s}^T (\Phi_{\{1\} \cup S, \{1\} \cup S})^{-1} e_{1,1+s} = \mathbf{K}_{11}^{-1}.$$

Therefore, to show $e_{1,p}^T \Phi^{-1} e_{1,p} \geq e_{1,1+s}^T (\Phi_{\{1\} \cup S, \{1\} \cup S})^{-1} e_{1,1+s}$, it suffices to show $\mathbf{K}_{11} \geq \mathbf{K}_{11} - \mathbf{K}_{21}^T \mathbf{K}_{22}^{-1} \mathbf{K}_{21}$. However, this is immediate to see since \mathbf{K}_{22} is positive definite. Besides, the equality holds if and only if $\mathbf{K}_{21} = 0$, or equivalently

$$\Phi_{S, \{1\}} = \Phi_{S, (\{1\} \cup S)^c} (\Phi_{(\{1\} \cup S)^c, (\{1\} \cup S)^c})^{-1} \Phi_{(\{1\} \cup S)^c, \{1\}}.$$

This completes the proof.

S5.3. Proof of Lemma S.2. To prove Lemma S.2, we need to show

$$(S5.30) \quad \Pr \left(\max_{1 \leq j \leq p} \|\mathbf{X}^j\|_\infty = O(\sqrt{n/\log(p)}) \right) \rightarrow 1,$$

$$(S5.31) \quad \Pr \left(\max_{1 \leq j \leq p} \|\mathbf{X}^j\|_2 = O(\sqrt{n}) \right) \rightarrow 1,$$

$$(S5.32) \quad \Pr (\lambda_{\min} (\mathbf{X}_{S \cup \mathcal{M}}^T \mathbf{X}_{S \cup \mathcal{M}}) \geq cn) \rightarrow 1,$$

$$(S5.33) \quad \Pr (\lambda_{\max} (\mathbf{X}_{S \cup \mathcal{M}}^T \mathbf{X}_{S \cup \mathcal{M}}) = O(n)) \rightarrow 1,$$

$$(S5.34) \quad \Pr (\|\mathbf{X}_{(S \cup \mathcal{M})^c}^T \mathbf{X}_{S \cup \mathcal{M}}\|_{2, \infty} = O(n)) \rightarrow 1,$$

$$(S5.35) \quad \Pr \left(\max_{1 \leq j \leq p} \lambda_{\max} \{ \mathbf{X}_{S \cup \mathcal{M}}^T \text{diag}(\|\mathbf{X}^j\|) \mathbf{X}_{S \cup \mathcal{M}} \} = O(n) \right) \rightarrow 1,$$

for some constant $c > 0$.

Under the given conditions, we have

$$\max_{1 \leq j \leq p} \|\mathbf{X}^j\|_\infty \leq \omega_0 \quad \text{and} \quad \max_{1 \leq j \leq p} \|\mathbf{X}^j\|_2 \leq \omega_0 \sqrt{n}.$$

Since $\log p = O(n^a)$ for some $0 < a < 1$, we have $\log p = o(n)$. This proves (S5.30) and (S5.31).

Note that

$$\begin{aligned} \frac{1}{n} \lambda_{\min}(\mathbf{X}_{S \cup \mathcal{M}}^T \mathbf{X}_{S \cup \mathcal{M}}) &= \frac{1}{n} \inf_{\substack{\mathbf{a} \in \mathbb{R}^{s+m} \\ \|\mathbf{a}\|_2=1}} \mathbf{a}^T \mathbf{X}_{S \cup \mathcal{M}}^T \mathbf{X}_{S \cup \mathcal{M}} \mathbf{a} \geq \inf_{\substack{\mathbf{a} \in \mathbb{R}^{s+m} \\ \|\mathbf{a}\|_2=1}} \mathbf{a}^T \mathbf{\Lambda}_{S \cup \mathcal{M}, S \cup \mathcal{M}} \mathbf{a} \\ &\quad - \left| \sup_{\substack{\mathbf{a} \in \mathbb{R}^{s+m} \\ \|\mathbf{a}\|_2=1}} \mathbf{a}^T \left(\frac{1}{n} \mathbf{X}_{i, S \cup \mathcal{M}} \mathbf{X}_{i, S \cup \mathcal{M}}^T - \mathbf{\Lambda}_{S \cup \mathcal{M}, S \cup \mathcal{M}} \right) \mathbf{a} \right| \\ &\geq \bar{c} - \left\| \frac{1}{n} \mathbf{X}_{i, S \cup \mathcal{M}} \mathbf{X}_{i, S \cup \mathcal{M}}^T - \mathbf{\Lambda}_{S \cup \mathcal{M}, S \cup \mathcal{M}} \right\|_2 \\ \text{(S5.36)} \quad &\geq \bar{c} - \left\| \frac{1}{n} \mathbf{X}_{i, S \cup \mathcal{M}} \mathbf{X}_{i, S \cup \mathcal{M}}^T - \mathbf{\Lambda}_{S \cup \mathcal{M}, S \cup \mathcal{M}} \right\|_\infty, \end{aligned}$$

where the third inequality follows by Lemma S.5 and the condition that $\lambda_{\min}(\mathbf{\Lambda}_{S \cup \mathcal{M}, S \cup \mathcal{M}}) \geq \bar{c}$, the last inequality is due to Lemma S.8.

For any $j_1, j_2 \in S \cup \mathcal{M}$, we have

$$\begin{aligned} &\mathbb{E} \exp(|X_{i, j_1} X_{i, j_2} - \mathbf{\Lambda}_{j_1, j_2}|) - 1 - |X_{i, j_1} X_{i, j_2} - \mathbf{\Lambda}_{j_1, j_2}| \\ &\leq \mathbb{E} \exp(|X_{i, j_1} X_{i, j_2} - \mathbf{\Lambda}_{j_1, j_2}|) \leq \mathbb{E} \exp(2|X_{i, j_1} X_{i, j_2}|) \leq \exp(2\omega_0^2), \end{aligned}$$

where the second inequality is due to Jensen's inequality. Similar to (S5.13), we can show

$$\begin{aligned} &\Pr \left(\max_{j_1, j_2 \in S \cup \mathcal{M}} \left| \frac{1}{n} \sum_{i=1}^n X_{i, j_1} X_{i, j_2} - \mathbf{\Lambda}_{j_1, j_2} \right| > \gamma \frac{\sqrt{\log n}}{\sqrt{n}} \right) \\ &\leq 2 \exp \left(-\frac{1}{2} \frac{\gamma^2 \log n}{\exp(2\omega_0^2) + \gamma \sqrt{n \log p}} + 2 \log(s+m) \right). \end{aligned}$$

Since $s+m = o(n)$, we can show for sufficiently large γ that

$$\text{(S5.37)} \quad \Pr \left(\max_{j_1, j_2 \in S \cup \mathcal{M}} \left| \frac{1}{n} \sum_{i=1}^n X_{i, j_1} X_{i, j_2} - \mathbf{\Lambda}_{j_1, j_2} \right| > \gamma \frac{\sqrt{\log n}}{\sqrt{n}} \right) \rightarrow 0.$$

This further implies

$$\Pr \left(\max_{j_1, j_2 \in \text{SUM}} \left\| \frac{1}{n} \sum_{i=1}^n X_{i, j_1} X_{i, j_2} - \mathbf{\Lambda}_{j_1, j_2} \right\|_{\infty} > \gamma(s+m) \frac{\sqrt{\log n}}{\sqrt{n}} \right) \rightarrow 0.$$

Under the given conditions, we have $(s+m)\sqrt{\log n} = o(\sqrt{n})$. Hence, we have for sufficiently large n ,

$$\Pr \left(\max_{j_1, j_2 \in \text{SUM}} \left\| \frac{1}{n} \sum_{i=1}^n X_{i, j_1} X_{i, j_2} - \mathbf{\Lambda}_{j_1, j_2} \right\|_{\infty} > \frac{\bar{c}}{2} \right) \rightarrow 0.$$

This together with (S5.36) implies that

$$(S5.38) \quad \Pr \left\{ \frac{1}{n} \lambda_{\min} (\mathbf{X}_{\text{SUM}}^T \mathbf{X}_{\text{SUM}}) > \frac{\bar{c}}{2} \right\} \rightarrow 1.$$

(S5.32) is hence proven. Similarly, we can show (S5.33) holds.

By condition, we have $\lambda_{\max}(\mathbf{\Lambda}_{\text{SUM}, \text{SUM}}) = O(1)$. It then follows from Cauchy-Schwarz inequality that

$$\begin{aligned} & \sup_{\substack{\mathbf{a} \in \mathbb{R}^{s+m} \\ \|\mathbf{a}\|_2 \leq 1}} |\mathbf{E} \mathbf{X}_{0, j} \mathbf{X}_{0, \text{SUM}}^T \mathbf{a}| \leq \sup_{\substack{\mathbf{a} \in \mathbb{R}^{s+m} \\ \|\mathbf{a}\|_2 \leq 1}} \sqrt{\mathbf{E} \mathbf{X}_{0, j}^2} \sqrt{\mathbf{E} (\mathbf{X}_{0, \text{SUM}}^T \mathbf{a})^2} \\ & \leq \omega_0^2 \sup_{\substack{\mathbf{a} \in \mathbb{R}^{s+m} \\ \|\mathbf{a}\|_2 \leq 1}} \mathbf{a}^T \mathbf{\Lambda}_{\text{SUM}, \text{SUM}} \mathbf{a} = O(\omega_0^2). \end{aligned}$$

This implies

$$(S5.39) \quad \left\| \mathbf{E} \mathbf{X}_{0, (\text{SUM})^c} \mathbf{X}_{0, \text{SUM}}^T \right\|_{2, \infty} = O(1).$$

Similar to (S5.37), we can show

$$\Pr \left(\max_{1 \leq j_1 \leq p, j_2 \in \text{SUM}} \left| \sum_{i=1}^n X_{i, j_1} X_{i, j_2} - n \mathbf{E} X_{0, j_1} X_{0, j_2} \right| > \gamma \sqrt{n \log p} \right) \rightarrow 0,$$

for some constant $\gamma > 0$. This further implies that

$$(S5.40) \quad \Pr \left(\left\| \sum_{i=1}^n \mathbf{X}_{i, (\text{SUM})^c} \mathbf{X}_{i, \text{SUM}} - n \mathbf{X}_{0, (\text{SUM})^c} \mathbf{X}_{0, \text{SUM}} \right\|_{2, \infty} > \gamma \sqrt{(s+m)n \log p} \right) \rightarrow 0.$$

Under the given conditions, we have $\sqrt{(s+m)n \log p} = O(n)$. This together with (S5.39) implies that

$$\Pr \left(\left\| \sum_{i=1}^n \mathbf{X}_{i,(S \cup \mathcal{M})^c} \mathbf{X}_{i,S \cup \mathcal{M}} \right\|_{2,\infty} = O(n) \right) \rightarrow 1.$$

(S5.34) is hence proven.

Finally, note that

$$\max_{1 \leq j \leq p} \lambda_{\max} (\mathbf{X}_{S \cup \mathcal{M}}^T \text{diag}\{|\mathbf{X}^j|\} \mathbf{X}_{S \cup \mathcal{M}}) \leq \omega_0 \lambda_{\max} (\mathbf{X}_{S \cup \mathcal{M}}^T \mathbf{X}_{S \cup \mathcal{M}}).$$

By (S5.33), we have

$$\Pr \left(\max_{1 \leq j \leq p} \lambda_{\max} (\mathbf{X}_{S \cup \mathcal{M}}^T \text{diag}\{|\mathbf{X}^j|\} \mathbf{X}_{S \cup \mathcal{M}}) = O(n) \right) \rightarrow 1.$$

This proves (S5.35). The proof is hence completed.

S5.4. Proof of Lemma S.3. It suffices to show (S5.30)-(S5.35) hold. By Bonferroni's inequality, Markov's inequality and the definition of the Orlicz norm, we have

$$\begin{aligned} & \Pr \left(\max_{1 \leq i \leq n, 1 \leq j \leq p} |X_{i,j}| > \sqrt{2}\omega_0 \sqrt{\log p + \log n} \right) \\ & \leq np \max_{1 \leq i \leq n, 1 \leq j \leq p} \Pr(|X_{i,j}| > \sqrt{2}\omega_0 \sqrt{\log p + \log n}) \\ & \leq np \mathbb{E} \frac{\exp(|X_{i,j}|^2/\omega_0^2)}{\exp(2\omega_0^2(\log p + \log n)/\omega_0^2)} \leq \frac{2np}{\exp(\log(np)^2)} = \frac{2}{np} \rightarrow 0. \end{aligned}$$

Therefore,

$$(S5.41) \quad \Pr \left(\max_{1 \leq i \leq n, 1 \leq j \leq p} |X_{i,j}| = O(\sqrt{\log p + \log n}) \right) \rightarrow 1.$$

Under the given conditions, we have $\log p + \log n = O(n/(\log p))$. This together with (S5.41) yields (S5.30).

By the definition of the Orlicz norm, we have

$$(S5.42) \quad \max_{1 \leq j \leq p} \mathbb{E} X_{0,j}^2 \leq \max_{1 \leq j \leq p} \|X_{0,j}^2\|_{\psi_1} = \max_{1 \leq j \leq p} \|X_{0,j}\|_{\psi_2}^2 \leq \omega_0^2.$$

Besides, it follows from the Bernstein's inequality (see Lemma G.3 in [Shi et al., 2017](#)) and Bonferroni's inequality that

$$\begin{aligned} & \Pr \left(\max_{j=1}^p \left| \sum_{i=1}^n X_{i,j}^2 - nEX_{0,j}^2 \right| \leq \gamma \sqrt{n \log p} \right) \\ & \leq \max_{j=1}^p p \Pr \left(\left| \sum_{i=1}^n X_{i,j}^2 - nEX_{0,j}^2 \right| \leq \gamma \sqrt{n \log p} \right) \rightarrow 1, \end{aligned}$$

for some constant $\gamma > 0$. Under the given conditions, this implies that

$$\Pr \left(\max_{j=1}^p \left| \sum_{i=1}^n X_{i,j}^2 - nEX_{0,j}^2 \right| \leq \gamma n \right) \rightarrow 1.$$

Combining this together with (S5.42) yields (S5.31).

It follows from Lemma C.1 in [Shi et al. \(2017\)](#) that we have with probability tending to 1,

$$(S5.43) \sup_{\substack{\mathbf{a} \in \mathbb{R}^{s+m} \\ \|\mathbf{a}\|_2=1}} \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{a}^T \mathbf{X}_{i,SUM})^2 - \mathbf{a}^T \mathbf{\Lambda}_{SUM,SUM} \mathbf{a} \right| = O \left(\frac{\sqrt{s \log n}}{\sqrt{n}} \right).$$

Note that the RHS of (S5.43) is $o(1)$. By condition, $\lambda_{\min}(\mathbf{\Lambda}_{SUM,SUM}) \geq \bar{c}$. Following the arguments in (S5.36), we have

$$\Pr \left(\lambda_{\min} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{X}_{i,SUM} \mathbf{X}_{i,SUM}^T \right) \geq \frac{\bar{c}}{2} \right) \rightarrow 1.$$

This proves (S5.32). Similarly, we can show (S5.33) holds.

Similar to (S5.42), we can show

$$\sup_{\mathbf{a} \in \mathbb{R}^{s+m}, \|\mathbf{a}\|_2 \leq 1} \mathbb{E} |\mathbf{a}^T \mathbf{X}_{0,SUM}|_2^2 \leq \sup_{\mathbf{a} \in \mathbb{R}^{s+m}, \|\mathbf{a}\|_2 \leq 1} \|\mathbf{a}^T \mathbf{X}_{0,SUM}\|_{\psi_2}^2 \leq \omega_0^2.$$

Hence, we have

$$\lambda_{\max}(\mathbf{\Lambda}_{SUM,SUM}) = O(1).$$

Using similar arguments in (S5.39) and (S5.40), we can show (S5.34) holds.

Finally, by (S5.42), we have

$$\max_{1 \leq j \leq p} \|X_{0,j}\|_{\psi_1} \leq \max_{1 \leq j \leq p} \|X_{0,j}\|_{\psi_2} / \sqrt{\log 2} \leq \omega_0 / \sqrt{\log 2}.$$

As a result, (S5.35) is directly implies by Lemma C.2 in [Shi et al. \(2017\)](#). This completes the proof.

APPENDIX S6: TECHNICAL LEMMAS

LEMMA S.4. Let $\widehat{S}_0 = \{j \in \mathcal{M}^c : \hat{\beta}_{0,j} \neq 0\}$. Suppose that $\hat{\beta}_0$ satisfies the following conditions. There exists some vector $\boldsymbol{\nu} \in \mathbb{R}^r$ such that

$$(S6.1) \quad \mathbf{X}_{\mathcal{M}}^T (\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}_{\widehat{S}_0 \cup \mathcal{M}} \hat{\beta}_{0, \widehat{S}_0 \cup \mathcal{M}})) = \mathbf{C}^T \boldsymbol{\nu},$$

$$(S6.2) \quad \mathbf{C} \hat{\beta}_{0, \mathcal{M}} = \mathbf{t},$$

$$(S6.3) \quad \mathbf{X}_{\widehat{S}}^T (\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X}_{\widehat{S}_0 \cup \mathcal{M}} \hat{\beta}_{0, \widehat{S}_0 \cup \mathcal{M}})) = n \bar{\rho}(\hat{\beta}_{0, \widehat{S}}, \lambda_{n,0}),$$

$$(S6.4) \quad \left\| \mathbf{X}_{(\widehat{S}_0 \cup \mathcal{M})^c}^T \{\mathbf{Y} - \boldsymbol{\mu}(\mathbf{X} \hat{\beta}_0)\} \right\|_{\infty} < \lambda_{n,0} \rho'(0+).$$

Furthermore, for any basis matrix $\mathbf{Z} \in \mathbb{R}^{m \times (m-r)}$ for the null space of \mathbf{C} ,

$$(S6.5) \quad \lambda_{\min} \left\{ \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_{\widehat{S}}^T \end{pmatrix} \boldsymbol{\Sigma}(\mathbf{X} \hat{\beta}_0) \begin{pmatrix} \mathbf{Z}^T \mathbf{X}_{\mathcal{M}}^T \\ \mathbf{X}_{\widehat{S}}^T \end{pmatrix}^T \right\} > n \lambda_{n,0} \kappa(\rho, \hat{\beta}_{0, \widehat{S}}, \lambda_{n,0}).$$

Then $\hat{\beta}_0$ is a local maximizer of $Q_n(\boldsymbol{\beta}, \lambda_{n,0})$ with the constraint $\mathbf{C} \boldsymbol{\beta}_{\mathcal{M}} = \mathbf{t}$.

REMARK S6.1. Note that solving (2.2) is equivalent to maximizing the following unconstrained problem:

$$\arg \max_{\boldsymbol{\beta}, \boldsymbol{\nu}} = \left(\frac{1}{n} \sum_{i=1}^n \{Y_i \boldsymbol{\beta}^T \mathbf{X}_i - b(\boldsymbol{\beta}^T \mathbf{X}_i)\} - \sum_{j \notin \mathcal{M}} p_{\lambda_{n,0}}(|\beta_j|) - \boldsymbol{\nu}^T (\mathbf{C} \boldsymbol{\beta}_{\mathcal{M}} - \mathbf{t}) \right).$$

Conditions (S6.1), (S6.2), (S6.3) and (S6.5) in Lemma S.4 guarantee that $\hat{\beta}_0$ is a local optimum to (2.2) when constrained on the subspace $\{\boldsymbol{\beta} : \boldsymbol{\beta}_{(\widehat{S}_0 \cup \mathcal{M})^c} = 0\}$ of \mathbb{R}^p . Condition (S6.4) ensures that $\hat{\beta}_0$ is indeed a local maximizer of (2.2) with the linear constraints.

Proof: We first show (S6.1), (S6.2), (S6.3) and (S6.5) guarantee that $\hat{\beta}_0$ is a local maximizer of $Q_n(\boldsymbol{\beta}, \lambda_{n,0})$ subject to $\mathbf{C} \boldsymbol{\beta}_{\mathcal{M}} = \mathbf{t}$ when constrained on the subspace $\{\boldsymbol{\beta} : \boldsymbol{\beta}_{(\widehat{S}_0 \cup \mathcal{M})^c} = 0\}$ of \mathbb{R}^p . For simplicity, we assume ρ is twice differentiable at $|\hat{\beta}_{0,j}|$ for any $j \in \mathcal{S}$. Since \mathbf{Z} is a basis matrix for the null space of \mathbf{C} , the matrix

$$\mathbf{L} = \begin{pmatrix} \mathbf{Z} & \mathbf{O}_{m \times s} \\ \mathbf{O}_{s \times m-r} & \mathbf{I}_s \end{pmatrix}$$

is a basis matrix for the null space of $(\mathbf{C} \mathbf{O}_{r \times s})$. The local optimality of

$\hat{\beta}_{0, \mathcal{M} \cup \hat{\mathcal{S}}}$ requires:

$$(S6.6) \quad \mathbf{C}\hat{\beta}_{0, \mathcal{M}} = \mathbf{t},$$

$$(S6.7) \quad \frac{\partial Q_n(\hat{\beta}_0, \lambda_{n,0})}{\partial \beta_{\hat{\mathcal{S}}}} = \mathbf{0},$$

$$(S6.8) \quad \frac{\partial Q_n(\hat{\beta}_0, \lambda_{n,0})}{\partial \beta_{\mathcal{M}}} = \mathbf{C}^T \boldsymbol{\nu},$$

$$(S6.9) \quad \lambda_{\min} \left\{ \mathbf{L}^T \begin{pmatrix} -\frac{\partial^2 Q_n(\hat{\beta}_0, \lambda_{n,0})}{\partial \beta_{\mathcal{M}} \partial \beta_{\mathcal{M}}^T} & -\frac{\partial^2 Q_n(\hat{\beta}_0, \lambda_{n,0})}{\partial \beta_{\mathcal{M}} \partial \beta_{\hat{\mathcal{S}}}^T} \\ -\frac{\partial^2 Q_n(\hat{\beta}_0, \lambda_{n,0})}{\partial \beta_{\hat{\mathcal{S}}} \partial \beta_{\mathcal{M}}^T} & -\frac{\partial^2 Q_n(\hat{\beta}_0, \lambda_{n,0})}{\partial \beta_{\hat{\mathcal{S}}} \partial \beta_{\hat{\mathcal{S}}}^T} \end{pmatrix} \mathbf{L} \right\} > 0.$$

It is immediate to see that, (S6.1), (S6.2) and (S6.3) directly imply (S6.6), (S6.7) and (S6.8). In the following, we show that (S6.5) implies (S6.9). When ρ is twice differentiable at $|\hat{\beta}_{0,j}|$ for any $j \in \mathcal{S}$, the matrix in the left-hand side (LHS) of (S6.9) is equal to

$$\begin{aligned} & \begin{pmatrix} \mathbf{Z}^T \mathbf{X}^T_{\mathcal{M}} \\ \mathbf{X}^T_{\hat{\mathcal{S}}} \end{pmatrix} \boldsymbol{\Sigma}(\mathbf{X}\hat{\beta}_0) \begin{pmatrix} \mathbf{Z}^T \mathbf{X}^T_{\mathcal{M}} \\ \mathbf{X}^T_{\hat{\mathcal{S}}} \end{pmatrix}^T \\ & - n\lambda_{n,0} \text{diag} \left(\underbrace{0, \dots, 0}_m, -\rho''(|\hat{\beta}_{0, \hat{\mathcal{S}}_1}|), \dots, -\rho''(|\hat{\beta}_{0, \hat{\mathcal{S}}_s}|) \right) \triangleq \mathbf{M}_1 - \mathbf{M}_2, \end{aligned}$$

where $\hat{\mathcal{S}}_j$ denotes the j th element in the set $\hat{\mathcal{S}}$. When ρ is not twice differentiable, we can replace \mathbf{M}_2 by a diagonal matrix whose absolute value is bounded by $\kappa(\rho, \hat{\beta}_{0, \hat{\mathcal{S}}}, \lambda_{n,0})$. Hence, LHS of (S6.9) is larger than $\lambda_{\min}(\mathbf{M}_1) - \lambda_{\max}(\mathbf{M}_2) \geq \lambda_{\min}(\mathbf{M}_1) - n\lambda_n \kappa(\rho, \hat{\beta}_{0, \hat{\mathcal{S}}}, \lambda_{n,0})$. Therefore, it follows from (S6.5) that (S6.9) is satisfied.

Under Condition (S6.4), using similar arguments in the proof of Theorem 1 in Lv and Fan (2009) or the proof of Theorem 1 in Fan and Lv (2011), we can show $\hat{\beta}_0$ is indeed a local maximizer of $Q_n(\beta, \lambda_{n,0})$ with $\mathbf{C}\beta_{\mathcal{M}} = \mathbf{t}$. This completes the proof.

LEMMA S.5. *For any symmetric matrix $\mathbf{A} \in \mathbb{R}^{q \times q}$, we have*

$$\|\mathbf{A}\|_2 = \sup_{\mathbf{a}: \|\mathbf{a}\|_2=1} |\mathbf{a}^T \mathbf{A} \mathbf{a}|.$$

Proof: By Cauchy-Schwarz inequality, we have

$$\sup_{\mathbf{a}: \|\mathbf{a}\|_2=1} |\mathbf{a}^T \mathbf{A} \mathbf{a}| \leq \sup_{\mathbf{a}: \|\mathbf{a}\|_2=1} \|\mathbf{a}\|_2^2 \|\mathbf{A}\|_2 = \|\mathbf{A}\|_2.$$

Hence, we have shown $\sup_{\mathbf{a}: \|\mathbf{a}\|_2=1} |\mathbf{a}^T \mathbf{A} \mathbf{a}| \leq \|\mathbf{A}\|_2$. It remains to show $\|\mathbf{A}\|_2 \leq \sup_{\mathbf{a}: \|\mathbf{a}\|_2=1} |\mathbf{a}^T \mathbf{A} \mathbf{a}|$. Since \mathbf{A} is symmetric, according to the eigen decomposition theorem, we have $\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$ for some orthogonal matrix \mathbf{U} and diagonal matrix $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_q)$. By definition, we have

$$(S6.10) \quad \|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A} \mathbf{A}^T)} = \sqrt{\lambda_{\max}(\mathbf{\Lambda}^2)} = \max_{j=1}^q |\lambda_j|.$$

Assume $\mathbf{U} = (u_1, \dots, u_q)$. Since \mathbf{U} is orthogonal, we have $u_i^T u_j = 0$ when $i \neq j$ and $u_j^T u_j = 1$. Hence, we have

$$\sup_{\mathbf{a}: \|\mathbf{a}\|_2=1} |\mathbf{a}^T \mathbf{A} \mathbf{a}| \geq \max_j |u_j^T \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T u_j| = \max_j |\lambda_j|.$$

This together with (S6.10) implies that $\sup_{\mathbf{a}: \|\mathbf{a}\|_2=1} |\mathbf{a}^T \mathbf{A} \mathbf{a}| \geq \|\mathbf{A}\|_2$. The proof is hence completed.

LEMMA S.6. *Let $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ be independent m -dimensional random vectors with $E\mathbf{X}_j = \mathbf{0}$, $\sum_j \text{cov}(\mathbf{X}_j) = \mathbf{I}_m$, let \mathbf{Z} denote an m -dimensional multivariate normal vector with mean $\mathbf{0}$ and covariance matrix \mathbf{I}_m , then*

$$\sup_C \left| \Pr\left(\sum \mathbf{X}_i \in C\right) - \Pr(\mathbf{Z} \in C) \right| \leq c_0 m^{1/4} \sum E\|\mathbf{X}_i\|_2^3,$$

for some constant c_0 , where the supremum is taken for all convex subsets in \mathbb{R}^m .

Proof: This follows directly from Theorem 1 in Bentkus (2004).

LEMMA S.7. *Denoted by $\chi^2(r, \gamma)$ a χ^2 random variable with r degrees of freedom and a non-centrality parameter γ . Then we have*

$$\lim_{\varepsilon \rightarrow 0^+} \sup_{r \geq 1, \gamma \geq 0} |\Pr(\chi^2(r, \gamma) \leq x + r\varepsilon) - \Pr(\chi^2(r, \gamma) \leq x - r\varepsilon)| \rightarrow 0.$$

Proof: It suffices to show

$$(S6.11) \quad \lim_{\varepsilon \rightarrow 0^+} \sup_{r \geq 1} |\Pr(\chi^2(r, 0) \leq x + r\varepsilon) - \Pr(\chi^2(r, 0) \leq x - r\varepsilon)| \rightarrow 0,$$

since

$$\begin{aligned}
& |\Pr(\chi^2(r, \gamma) \leq x + r\varepsilon) - \Pr(\chi^2(r, \gamma) \leq x - r\varepsilon)| \\
& \leq \left| \sum_{k=0}^{\infty} \exp(-\gamma/2) \frac{(\gamma/2)^k}{k!} \{ \Pr(\chi^2(r + 2k, 0) \leq x + r\varepsilon) - \Pr(\chi^2(r + 2k, 0) \leq x - r\varepsilon) \} \right| \\
& \leq \sum_{k=0}^{\infty} \exp(-\gamma/2) \frac{(\gamma/2)^k}{k!} |\Pr(\chi^2(r + 2k, 0) \leq x + r\varepsilon) - \Pr(\chi^2(r + 2k, 0) \leq x - r\varepsilon)| \\
& \leq \sup_r |\Pr(\chi^2(r, 0) \leq x + r\varepsilon) - \Pr(\chi^2(r, 0) \leq x - r\varepsilon)|.
\end{aligned}$$

Below, we show (S6.11) holds. We first prove

$$(S6.12) \lim_{\varepsilon \rightarrow 0^+} \sup_{r \geq 2, x} |\Pr(\chi^2(r, 0) \leq x + r\varepsilon) - \Pr(\chi^2(r, 0) \leq x - r\varepsilon)| \rightarrow 0.$$

Note that (S6.12) holds when the probability density function f_r of a χ^2 distribution with r degrees of freedom satisfies $\sup_{r \geq 2} \sup_x r f_r(x) = O(1)$. By definition, we have

$$f_r(x) = \frac{1}{2^{r/2} \Gamma(r/2)} x^{r/2-1} \exp(-x/2).$$

The supremum of f_r is achieved at $x = r - 2$. For $r = 2$, obviously we have $\sup_x 2|f_2(x)| \leq 2|f_2(0)| = 1$. For $r \geq 3$, we have

$$(S6.13) \sup_{x, r \geq 3} |f_r(x)| \leq \sup_{r \geq 3} f_r(r - 2) \leq \sup_{r \geq 3} \frac{(r/2 - 1)^{r/2-1}}{\Gamma(r/2)} \exp(-r/2 + 1).$$

By Stirling's formula, we have

$$\Gamma(r/2) \geq \sqrt{2\pi} \left(\frac{r}{2} - 1\right)^{r/2-1/2} \exp\left(-\frac{r}{2} + 1\right),$$

which together with (S6.13) implies that

$$\sup_x \sup_{r \geq 3} r |f_r(x)| \leq \frac{r}{\sqrt{2\pi}(r/2 - 1)} \leq \frac{3\sqrt{2}}{\sqrt{\pi}}.$$

This proves (S6.12). It remains to show

$$(S6.14) \lim_{\varepsilon \rightarrow 0^+} \sup_x |\Pr(\chi^2(1, 0) \leq x + \varepsilon) - \Pr(\chi^2(1, 0) \leq x - \varepsilon)| \rightarrow 0.$$

For any $\epsilon > 0$, there exists some $\delta > 0$ such that $\Pr(\chi^2(1, 0) \leq 2\delta) \leq \epsilon$. Hence, for any $x \leq \delta$ and any $\epsilon \leq \delta$, we have

$$(S6.15) \quad \begin{aligned} & \Pr(\chi^2(1, 0) \leq x + \epsilon) - \Pr(\chi^2(1, 0) \leq x - \epsilon) \\ & \leq \Pr(\chi^2(1, 0) \leq 2\delta) \leq \epsilon. \end{aligned}$$

Moreover, $f_1(x)$ is continuous and monotonically decreasing on $[\delta, +\infty)$. This implies $\sup_{x \geq \delta} f_1(x) \leq f_1(\delta)$ and hence

$$\sup_{x \geq \delta} |\Pr(\chi^2(1, 0) \leq x + \epsilon) - \Pr(\chi^2(1, 0) \leq x - \epsilon)| \leq \frac{2\epsilon}{f_1(\delta)}.$$

Therefore, for any $x \geq \delta$ and any $0 < \epsilon \leq (2\epsilon)/\{f_1(\delta)\}$, we have

$$(S6.16) \quad \begin{aligned} & \Pr(\chi^2(1, 0) \leq x + \epsilon) - \Pr(\chi^2(1, 0) \\ & \leq x - \epsilon) \leq \Pr(\chi^2(1, 0) \leq 2\delta) \leq \epsilon. \end{aligned}$$

Equation (S6.14) now follows from (S6.15) and (S6.16). This completes the proof.

LEMMA S.8. *For any symmetric matrix \mathbf{A} , we have $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_\infty$.*

Proof : Note that $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_\infty \|\mathbf{A}\|_1$. Since \mathbf{A} is symmetric, we have $\|\mathbf{A}\|_\infty = \|\mathbf{A}\|_1$. This proves the assertion.

LEMMA S.9 (Matrix inversion in block form). *For any positive definite matrix*

$$\Psi = \begin{pmatrix} \Psi_{11} & \Psi_{12} \\ \Psi_{21} & \Psi_{22} \end{pmatrix},$$

denote its inverse matrix as Ω and partition it into $\Omega_{11}, \dots, \Omega_{22}$ accordingly. Then

$$\Omega_{11} = (\Psi_{11} - \Psi_{12}\Psi_{22}^{-1}\Psi_{21})^{-1}, \quad \Omega_{22} = (\Psi_{22} - \Psi_{21}\Psi_{11}^{-1}\Psi_{12})^{-1}.$$

As a result, the matrices $\Omega_{11} - \Psi_{11}^{-1}$ and $\Omega_{22} - \Psi_{22}^{-1}$ are positive semidefinite.

APPENDIX S7: REAL DATA ANALYSIS

In this section, we apply our proposed testing procedures to the European American single nucleotide polymorphisms (SNPs) data set (Price et al., 2006), which consists of 488 European American samples. We use the height phenotype (0/1, binary variable) of these European American samples as

the response, and focus on finding variables that are associated with this phenotype among a set of 277 SNPs. The genotype for each SNP is a categorical variable, coded as 0/1/2. We removed the outlier individuals as in Price et al. (2006). This gives us a total of 361 observations. However, for each observation, approximately 2% of SNPs are missing on average. We imputed all the missing values using the R package `missForest` available in CRAN. This package uses a random forest trained based on the observed entries to predict those missing values.

To formulate the testing hypotheses, we adopt a data splitting procedure. More specifically, we first randomly sample 20% of the observations and perform some preliminary analysis based on these observations. Results are given in Section S7.1 and we find 13 SNPs that are highly correlated with the phenotype. Based on the remaining 80% of the observations, we focus on testing whether the regression coefficients of these 13 variables are zero using logistic regression. The p-values of the partial penalized Wald, score and likelihood ratio statistics are 6.8×10^{-3} , 7.4×10^{-3} and 7.0×10^{-3} respectively. Under the significance level of 0.05, we reject the null hypothesis and conclude that at least one of the 13 regression coefficients is not equal to zero.

Since each covariate X_i^j is discrete and takes value on $\{0, 1, 2\}$, we can define the dummy variables $Z_i^{(j,1)} = I(X_i^j = 1)$, $Z_i^{(j,2)} = I(X_i^j = 2)$ and use these $\mathbf{Z}^{(j,m)}$'s as covariates in our logistic regression model. This yields a total of 554 covariates. Denoted by $\beta_{j,m}$ the corresponding regression coefficient of $\mathbf{Z}^{(j,m)}$, our next goal is to test the following hypothesis based on the remaining 80% of the samples:

$$H_0 : 2\beta_{j,1} = \beta_{j,2}, \quad \forall j \in \mathcal{M},$$

where \mathcal{M} denotes the set of the 13 SNPs selected in the preliminary analysis. Under H_0 , we have $\beta_{j,1}\mathbf{X}^j = \beta_{j,1}\mathbf{Z}^{(j,1)} + \beta_{j,2}\mathbf{Z}^{(j,2)}$ for any $j \in \mathcal{M}$. This corresponds to a lack of fit test regarding these important variables. The p-values of the Wald, score and likelihood ratio statistics are 0.083, 0.069 and 0.086 respectively. Hence, we fail to reject H_0 .

S7.1. Preliminary analysis in the real data application. We apply sure independence screening (Fan and Song, 2010) across all 277 variables. Specifically, we independently fit 277 logistic regressions by maximizing the marginal likelihood with the response and each univariate covariate and obtain the p-values from each marginal model. We find out there are a total of 13 variables with p-values smaller than 0.05. We report these variables and their associated p-values in the table below.

SNP	22	50	68	126	134	138	141	151	152	187	199	260	272
p-value (%)	0.4	3.5	4.2	1.9	4.0	4.8	4.3	4.6	3.5	4.1	4.5	3.9	1.8

APPENDIX S8: ADDITIONAL SIMULATION RESULTS

S8.1. Simulations results for Poisson regression. We consider Poisson regression in this section. The data were generated from the following model,

$$\Pr(Y_i = k | \mathbf{X}_i) = \frac{\exp(-\lambda_i) \lambda_i^k}{k!},$$

where

$$\lambda_i = \exp\left(0.75X_i^{(1)} - (0.75 + h^{(1)})X_i^{(2)} + h^{(2)}(X_i^{(3)} + X_i^{(4)})\right),$$

and $\mathbf{X}_i \stackrel{iid}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma})$ for some $p \times p$ covariance matrix $\boldsymbol{\Sigma}$. Consider testing the following hypotheses:

$$\begin{aligned} H_0^{(1)} : \beta_{1,0} + \beta_{2,0} = 0, \quad v.s \quad H_a^{(1)} : \beta_{1,0} + \beta_{2,0} \neq 0 \\ H_0^{(2)} : \beta_{3,0} = \beta_{4,0} = 0, \quad v.s \quad H_a^{(2)} : \beta_{3,0} \neq 0 \quad \text{or} \quad \beta_{4,0} \neq 0. \end{aligned}$$

We set $h^{(2)} = 0$ when testing $H_0^{(1)}$, and set $h^{(1)} = 0$ when testing $H_0^{(2)}$. We use the following four settings: (i) $p = 50$, $\boldsymbol{\Sigma} = \mathbf{I}_p$; (ii) $p = 200$, $\boldsymbol{\Sigma} = \mathbf{I}_p$; (iii) $p = 50$, $\boldsymbol{\Sigma} = \{0.5^{|i-j|}\}$; (iiii) $p = 200$, $\boldsymbol{\Sigma} = \{0.5^{|i-j|}\}$. For each setting, we set $h^{(j)} = 0, 0.03, 0.06, 0.12$ when testing $H_0^{(j)}$. The sample size is set to be 500. We provide the rejection probabilities for $H_0^{(1)}$ and $H_0^{(2)}$ in Table S7. The kernel density estimates of three test statistics under $H_0^{(1)}$ and $H_0^{(2)}$ are plotted in Figure S7 and Figure S8.

It can be seen that the Type-I error rates are well controlled under the null hypotheses, and the powers increase as $h^{(1)}$ or $h^{(2)}$ increases.

S8.2. Tables and plots.

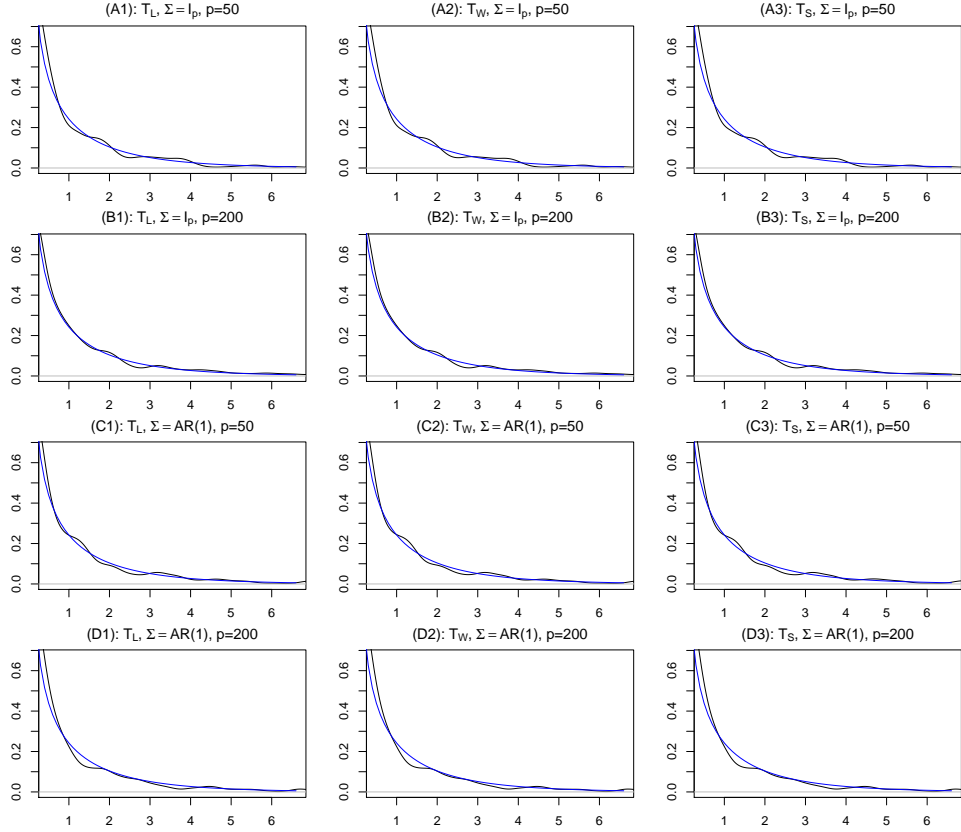


FIG S1. Kernel density plot of three test statistics under $H_0^{(1)}$ with different combinations of p and the covariance matrices. T_L , T_W and T_S from left to right. $p = 50, \Sigma = I_p$ and $p = 200, \Sigma = I_p$ and $p = 50, \Sigma = \{0.5^{|i-j|}\}_{i,j}$ and $p = 200, \Sigma = \{0.5^{|i-j|}\}_{i,j}$ from top to bottom. The black line plot the density function of a χ^2 distribution with 1 degree of freedom.

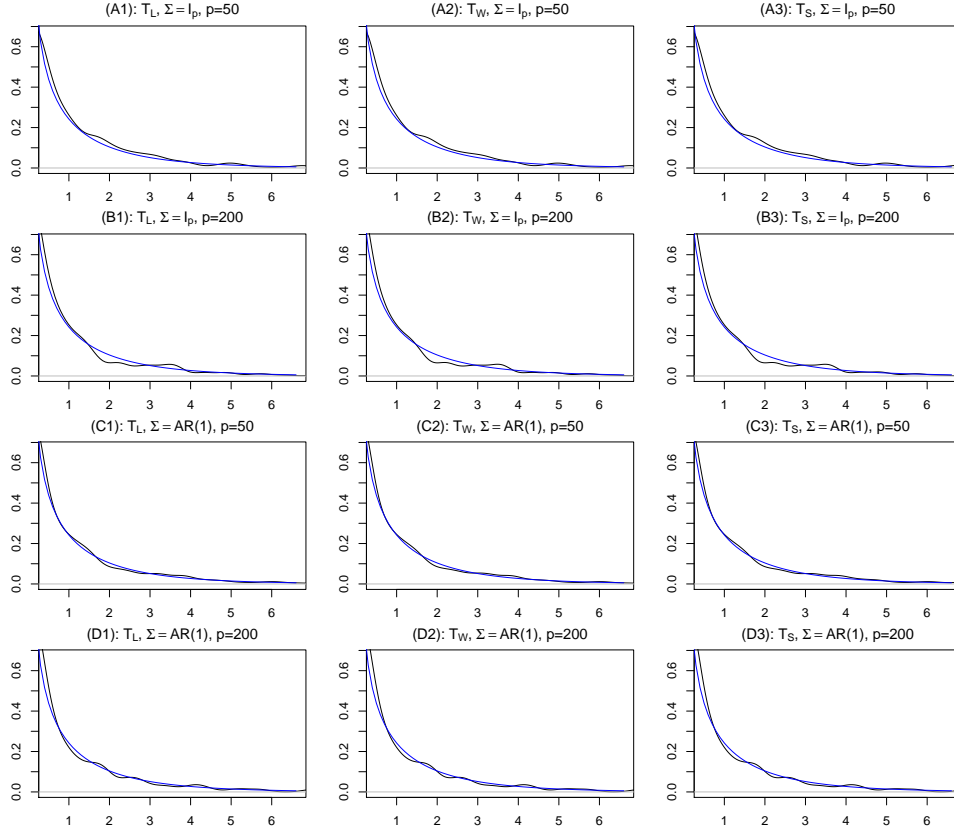


FIG S2. Kernel density plot of three test statistics under $H_0^{(2)}$ with different combinations of p and the covariance matrices. T_L , T_W and T_S from left to right. $p = 50, \Sigma = I_p$ and $p = 200, \Sigma = I_p$ and $p = 50, \Sigma = \{0.5^{|i-j|}\}_{i,j}$ and $p = 200, \Sigma = \{0.5^{|i-j|}\}_{i,j}$ from top to bottom. The black line plot the density function of a χ^2 distribution with 1 degree of freedom.

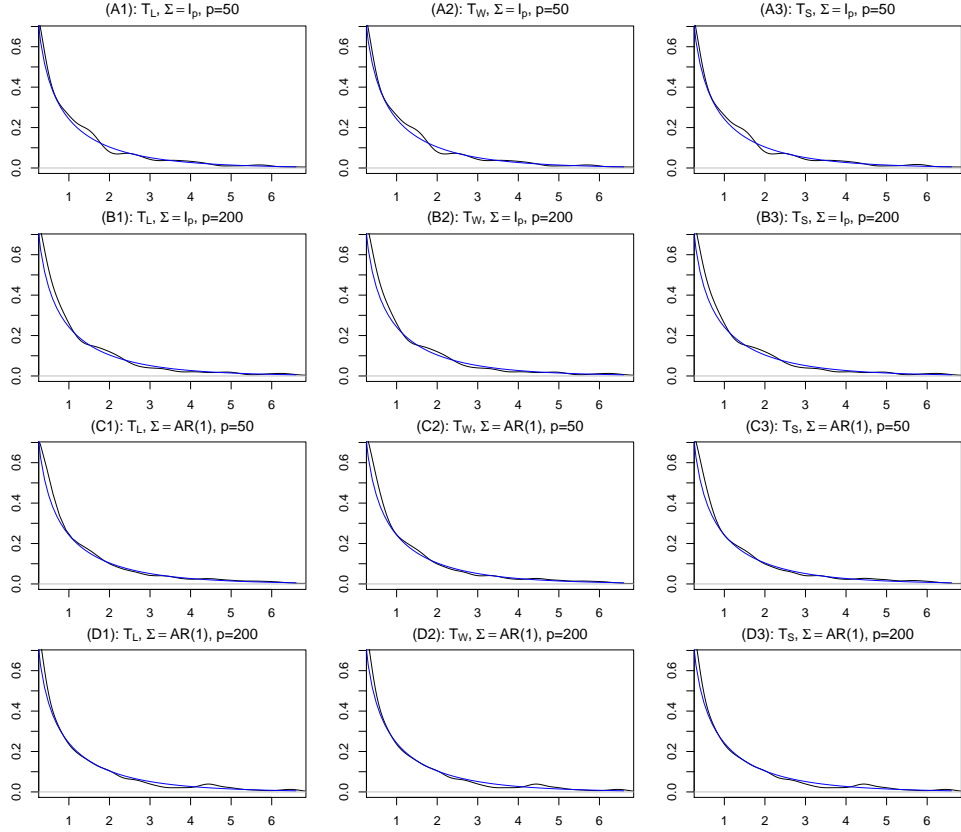


FIG S3. Kernel density plot of three test statistics under $H_0^{(2)}$ with different combinations of p and the covariance matrices. T_L , T_W and T_S from left to right. $p = 50, \Sigma = I_p$ and $p = 200, \Sigma = I_p$ and $p = 50, \Sigma = \{0.5^{|i-j|}\}_{i,j}$ and $p = 200, \Sigma = \{0.5^{|i-j|}\}_{i,j}$ from top to bottom. The black line plot the density function of a χ^2 distribution with 1 degree of freedom.

TABLE S1

Rejection probabilities (%) of the partial penalized Wald, score and likelihood ratio statistics with standard errors in parenthesis (%), under the setting where $\Sigma = \mathbf{I}$.

	$p = 50$			$p = 200$		
	T_L	T_W	T_S	T_L	T_W	T_S
$h^{(1)}$	$H_0^{(1)}$					
0	3.50(0.75)	3.50(0.75)	3.50(0.75)	6.00(0.97)	6.00(0.97)	6.00(0.97)
0.1	10.33(1.24)	10.33(1.24)	10.33(1.24)	9.17(1.18)	9.17(1.18)	9.17(1.18)
0.2	26.67(1.81)	26.67(1.81)	26.67(1.81)	32.00(1.90)	32.00(1.90)	32.00(1.90)
0.4	83.17(1.53)	83.17(1.53)	83.33(1.52)	78.83(1.67)	78.83(1.67)	78.83(1.67)
$h^{(2)}$	$H_0^{(2)}$					
0	4.67(0.86)	4.67(0.86)	4.83(0.88)	4.17(0.82)	4.17(0.82)	4.33(0.83)
0.1	13.50(1.40)	13.50(1.40)	13.67(1.40)	10.17(1.23)	10.17(1.23)	10.33(1.24)
0.2	32.83(1.92)	32.67(1.91)	32.83(1.92)	27.80(1.83)	28.80(1.83)	28.80(1.83)
0.4	80.50(1.62)	80.50(1.62)	80.50(1.62)	80.50(1.62)	80.67(1.61)	80.67(1.61)
$h^{(2)}$	$H_0^{(3)}$					
0	5.00(0.89)	5.00(0.89)	5.17(0.90)	5.17(0.90)	5.17(0.90)	5.17(0.90)
0.1	10.83(1.27)	10.83(1.27)	11.00(1.28)	11.33(1.29)	11.33(1.29)	11.50(1.30)
0.2	28.67(1.85)	28.67(1.85)	28.83(1.85)	32.50(1.91)	32.50(1.91)	32.50(1.91)
0.4	80.67(1.61)	80.67(1.61)	80.67(1.61)	81.50(1.59)	81.50(1.59)	81.50(1.59)

TABLE S2

Rejection probabilities (%) of the partial penalized Wald, score and likelihood ratio statistics, the Wald test statistic based on the de-sparsified Lasso estimator and the decorrelated score statistic under the settings where $\Sigma = \mathbf{I}$, with standard errors in parenthesis (%).

	T_L	T_W	T_S	T_W^D	T_S^D
$h^{(1)}$	$H_0^{(4)}$ and $p = 50$				
0	6.33(0.99)	6.50(1.01)	6.67(1.02)	5.17(0.90)	6.50(1.01)
0.1	18.67(1.59)	18.67(1.59)	18.83(1.60)	14.00(1.42)	18.50(1.59)
0.2	50.83(2.04)	50.83(2.04)	51.00(2.04)	41.33(2.01)	51.83(2.04)
0.4	95.67(0.83)	95.67(0.83)	95.67(0.83)	93.83(0.98)	96.00(0.80)
	$H_0^{(4)}$ and $p = 200$				
0	5.67(0.94)	5.67(0.94)	5.83(0.96)	5.17(0.90)	7.33(1.06)
0.1	15.67(1.48)	15.67(1.48)	15.67(1.48)	14.00(1.42)	15.67(1.48)
0.2	49.33(2.04)	49.33(2.04)	49.33(2.04)	41.33(2.01)	51.00(2.04)
0.4	97.33(0.66)	97.33(0.66)	97.33(0.66)	93.83(0.98)	97.67(0.62)
$h^{(2)}$	$H_0^{(5)}$ and $p = 50$				
0	6.33(0.99)	6.33(0.99)	6.50(1.00)	5.83(0.96)	6.50(1.01)
0.1	17.83(1.56)	17.83(1.56)	18.00(1.57)	14.17(1.42)	16.83(1.53)
0.2	53.00(2.04)	53.00(2.04)	53.00(2.04)	43.33(2.02)	52.17(2.04)
0.4	97.33(0.66)	97.33(0.66)	97.33(0.66)	95.17(0.88)	96.67(0.73)
	$H_0^{(5)}$ and $p = 200$				
0	4.00(0.80)	4.00(0.80)	4.00(0.80)	4.67(0.86)	3.17(0.71)
0.1	16.67(1.52)	16.67(1.52)	16.83(1.53)	13.67(1.40)	14.00(1.42)
0.2	49.17(2.04)	49.17(2.04)	49.17(2.04)	39.67(2.00)	43.17(2.02)
0.4	97.50(0.64)	97.50(0.64)	97.50(0.64)	92.50(1.08)	95.83(0.82)

TABLE S3

Rejection probabilities (%) of the partial penalized Wald, score and likelihood ratio statistics with standard errors in parenthesis (%), under the settings where $\Sigma = \mathbf{I}$.

	$p = 50$			$p = 200$		
	T_L	T_W	T_S	T_L	T_W	T_S
$h^{(1)}$	$H_0^{(6)}$					
0	5.50(0.93)	5.50(0.93)	5.67(0.93)	6.50(1.01)	6.50(1.01)	6.50(1.01)
0.2	15.00(1.46)	15.00(1.46)	15.67(1.46)	16.17(1.50)	16.17(1.50)	16.17(1.50)
0.4	50.33(2.04)	50.33(2.04)	50.50(2.04)	49.17(2.04)	49.17(2.04)	49.17(2.04)
0.8	97.50(0.64)	97.50(0.64)	97.50(0.64)	96.50(0.75)	96.50(0.75)	96.50(0.75)
$h^{(1)}$	$H_0^{(7)}$					
0	5.00(0.89)	5.00(0.89)	5.17(0.90)	5.83(0.96)	5.83(0.96)	5.83(0.96)
0.2	9.33(1.19)	9.17(1.18)	9.50(1.20)	13.00(1.37)	13.00(1.37)	13.00(1.37)
0.4	30.00(1.87)	30.00(1.87)	30.17(1.87)	26.83(1.81)	26.83(1.81)	26.83(1.81)
0.8	76.33(1.74)	76.67(1.73)	76.67(1.73)	76.67(1.73)	76.67(1.73)	76.67(1.73)
$h^{(1)}$	$H_0^{(8)}$					
0	4.83(0.88)	4.83(0.88)	5.00(0.89)	6.83(1.03)	6.83(1.03)	6.83(1.03)
0.2	8.83(1.16)	8.67(1.15)	9.00(1.17)	10.33(1.24)	10.33(1.24)	10.33(1.24)
0.4	21.33(1.67)	21.33(1.67)	21.50(1.68)	20.50(1.65)	20.50(1.65)	20.50(1.65)
0.8	55.83(2.03)	56.17(2.03)	56.00(2.03)	58.33(2.01)	58.33(2.01)	58.33(2.01)

TABLE S4
Rejection probabilities (%) of the partial penalized Wald, score and likelihood ratio statistics with standard errors in parenthesis (%).

	$p = 50$			$p = 200$		
	T_L	T_W	T_S	T_L	T_W	T_S
$h^{(1)}$	$H_0^{(1)}$ and $\Sigma = \mathbf{I}$					
0	5.33(0.92)	5.17(0.90)	5.33(0.92)	6.00(0.97)	5.33(0.92)	6.00(0.97)
0.2	11.33(1.29)	10.17(1.23)	10.83(1.27)	13.67(1.40)	13.00(1.37)	13.50(1.40)
0.4	33.33(1.92)	32.00(1.90)	33.17(1.92)	38.83(1.99)	37.67(1.98)	38.67(1.99)
0.8	88.50(1.30)	88.67(1.29)	88.83(1.29)	88.83(1.29)	88.17(1.32)	88.67(1.29)
	$H_0^{(1)}$ and $\Sigma = \{0.5^{ i-j }\}$					
0	6.50(1.01)	6.00(0.97)	6.33(0.99)	4.50(0.85)	4.33(0.83)	4.33(0.83)
0.2	19.17(1.61)	18.67(1.59)	19.00(1.60)	21.00(1.67)	20.33(1.64)	20.67(1.65)
0.4	61.50(1.99)	61.17(1.99)	61.50(1.99)	62.00(1.98)	62.17(1.98)	62.00(1.98)
0.8	98.33(0.52)	99.67(0.24)	99.67(0.24)	98.50(0.50)	99.33(0.33)	99.33(0.33)
$h^{(2)}$	$H_0^{(2)}$ and $\Sigma = \mathbf{I}$					
0	6.17(0.98)	6.00(0.97)	5.83(0.96)	5.67(0.94)	5.50(0.93)	5.83(0.96)
0.2	7.83(1.10)	9.17(1.12)	9.17(1.12)	8.00(1.11)	10.17(1.23)	10.00(1.22)
0.4	23.50(1.73)	26.50(1.80)	26.67(1.81)	23.83(1.74)	26.67(1.81)	26.83(1.81)
0.8	71.00(1.85)	76.00(1.74)	75.50(1.76)	74.00(1.79)	77.67(1.70)	77.50(1.70)
	$H_0^{(2)}$ and $\Sigma = \{0.5^{ i-j }\}$					
0	5.67(0.94)	5.33(0.92)	5.67(0.94)	4.50(0.85)	4.17(0.82)	4.17(0.82)
0.2	11.17(1.29)	12.67(1.36)	12.67(1.36)	11.00(1.28)	12.50(1.35)	12.50(1.35)
0.4	31.17(1.89)	35.33(1.95)	35.17(1.95)	33.00(1.92)	35.83(1.96)	35.67(1.96)
0.8	84.50(1.48)	88.33(1.31)	88.33(1.33)	88.67(1.29)	90.17(1.22)	90.00(1.22)
$h^{(2)}$	$H_0^{(3)}$ and $\Sigma = \mathbf{I}$					
0	4.17(0.82)	3.67(0.77)	3.67(0.77)	6.17(0.98)	5.67(0.94)	5.83(0.96)
0.2	8.67(1.15)	8.83(1.13)	8.83(1.13)	14.50(1.44)	13.67(1.40)	14.17(1.42)
0.4	35.83(1.96)	34.17(1.94)	35.50(1.95)	39.50(2.00)	39.00(1.99)	39.17(1.99)
0.8	90.00(1.22)	89.00(1.28)	89.67(1.24)	90.17(1.22)	89.67(1.24)	90.00(1.22)
	$H_0^{(3)}$ and $\Sigma = \{0.5^{ i-j }\}$					
0	4.83(0.88)	4.67(0.86)	4.67(0.86)	5.83(0.96)	5.33(0.92)	5.67(0.94)
0.2	18.00(1.57)	17.00(1.53)	17.67(1.56)	18.50(1.59)	18.17(1.57)	18.33(1.58)
0.4	55.17(2.03)	54.50(2.03)	55.00(2.03)	53.83(2.04)	53.33(2.04)	53.83(2.04)
0.8	98.33(0.52)	99.00(0.41)	99.00(0.41)	98.50(0.50)	98.33(0.44)	98.33(0.44)

TABLE S5

Rejection probabilities (%) of the partial penalized Wald, score and likelihood ratio statistics, the Wald test statistic based on the de-sparsified Lasso estimator and the decorrelated score statistic, with standard errors in parenthesis (%).

	T_L	T_W	T_S	T_W^D	T_S^D
$h^{(1)}$	$H_0^{(4)}, \Sigma = \mathbf{I}$ and $p = 50$				
0	4.83(0.86)	4.67(0.86)	4.67(0.86)	51.17(2.04)	7.33 (1.06)
0.2	12.50(1.35)	8.00(1.11)	8.33(1.12)	23.33(1.73)	18.67 (1.59)
0.4	27.17(1.82)	20.83(1.66)	21.33(1.67)	9.00(1.17)	37.17 (1.97)
0.8	77.00(1.72)	70.67(1.86)	71.83(1.84)	15.50(1.48)	82.67 (1.55)
	$H_0^{(4)}, \Sigma = \mathbf{I}$ and $p = 200$				
0	4.17(0.82)	4.33(0.83)	4.33(0.83)	70.67(1.86)	8.83(1.16)
0.2	12.00(1.33)	9.17(1.18)	9.17(1.18)	43.17(2.02)	22.33(1.70)
0.4	31.17(1.89)	23.33(1.73)	24.33(1.75)	24.00(1.74)	46.67(2.04)
0.8	78.33(1.67)	70.83(1.86)	71.83(1.84)	15.00(1.46)	89.33(1.26)
	$H_0^{(4)}, \Sigma = \{0.5^{ i-j }\}$ and $p = 50$				
0	6.00(0.97)	5.33(0.92)	5.50(0.93)	31.00(1.89)	7.00(1.04)
0.2	13.50(1.40)	10.83(1.27)	11.17(1.29)	11.00(1.28)	14.67(1.44)
0.4	33.83(1.93)	28.17(1.84)	29.17(1.86)	8.00(1.11)	38.83(1.99)
0.8	82.50(1.55)	78.33(1.68)	78.67(1.67)	37.33(1.97)	89.33(1.26)
	$H_0^{(4)}, \Sigma = \{0.5^{ i-j }\}$ and $p = 200$				
0	4.33(0.83)	4.00(0.80)	4.00(0.80)	57.83(2.02)	7.50(1.08)
0.2	10.67(1.26)	8.00(1.11)	8.50(1.14)	27.83(1.83)	18.67(1.59)
0.4	32.00(1.90)	27.00(1.81)	27.33(1.82)	13.17(1.38)	47.50(2.04)
0.8	79.67(1.64)	75.17(1.76)	75.17(1.76)	23.67(1.73)	88.50(1.30)
$h^{(2)}$	$H_0^{(5)}, \Sigma = \mathbf{I}$ and $p = 50$				
0	5.33(0.92)	5.00(0.89)	5.17(0.90)	2.50(0.64)	4.17(0.82)
0.2	19.83(1.63)	19.33(1.61)	19.83(1.63)	9.83(1.22)	17.67(1.56)
0.4	62.50(1.98)	60.67(1.99)	61.83(1.98)	42.00(2.01)	60.33(2.00)
0.8	99.50(0.29)	99.50(0.29)	99.50(0.29)	97.50(0.64)	99.33(0.33)
	$H_0^{(5)}, \Sigma = \mathbf{I}$ and $p = 200$				
0	5.00(0.89)	4.50(0.85)	4.83(0.88)	3.00(0.70)	3.00(0.70)
0.2	22.50(1.70)	21.83(1.69)	22.33(1.70)	13.17(1.38)	18.67(1.59)
0.4	64.00(1.96)	62.67(1.97)	63.67(1.96)	47.67(2.04)	58.17(2.01)
0.8	99.50(0.29)	99.67(0.24)	99.67(0.24)	98.00(0.57)	99.17(0.37)
	$H_0^{(5)}, \Sigma = \{0.5^{ i-j }\}$ and $p = 50$				
0	5.50(0.93)	4.83(0.88)	5.50(0.93)	3.67(0.77)	3.00(0.70)
0.2	21.33(1.67)	20.83(1.66)	21.00(1.66)	8.00(1.11)	15.33(1.47)
0.4	63.33(1.97)	62.67(1.97)	63.50(1.97)	34.67(1.94)	53.00(2.04)
0.8	99.17(0.37)	99.83(0.17)	99.83(0.17)	95.17(0.88)	99.17(0.37)
	$H_0^{(5)}, \Sigma = \{0.5^{ i-j }\}$ and $p = 200$				
0	7.67(1.09)	7.17(1.05)	7.67(1.09)	5.50(0.93)	4.00(0.80)
0.2	22.00(1.66)	21.00(1.69)	21.83(1.66)	9.00(1.17)	15.50(1.48)
0.4	64.50(1.95)	64.00(1.96)	64.33(1.96)	31.67(1.90)	48.00(2.04)
0.8	99.00(0.41)	99.50(0.29)	99.50(0.29)	94.33(0.94)	97.50(0.64)

TABLE S6
Rejection probabilities (%) of the partial penalized Wald, score and likelihood ratio statistics with standard errors in parenthesis (%).

	$p = 50$			$p = 200$		
	T_L	T_W	T_S	T_L	T_W	T_S
$h^{(1)}$	$H_0^{(6)}$ and $\Sigma = \mathbf{I}$					
0	5.67(0.94)	5.17(0.90)	5.50(0.93)	5.17(0.90)	4.83(0.88)	5.17(0.90)
0.4	19.17(1.61)	18.67(1.59)	19.17(1.61)	21.50(1.68)	21.00(1.66)	21.33(1.67)
0.8	61.50(1.99)	61.17(1.99)	61.33(1.99)	59.00(2.01)	57.50(2.02)	58.83(2.01)
1.6	97.67(0.62)	97.83(0.59)	97.83(0.59)	98.00(0.57)	98.00(0.57)	98.17(0.55)
	$H_0^{(6)}$ and $\Sigma = \{0.5^{ i-j }\}$					
0	4.33(0.83)	3.83(0.78)	4.33(0.80)	4.83(0.88)	4.67(0.86)	4.83(0.88)
0.4	45.33(2.03)	44.67(2.03)	45.17(2.03)	45.50(2.03)	44.33(2.03)	45.00(2.03)
0.8	93.50(1.01)	94.00(0.97)	94.17(0.96)	93.00(1.04)	92.67(1.06)	93.00(1.04)
1.6	100.00(0.00)	100.00(0.00)	100.00(0.00)	100.00(0.00)	100.00(0.00)	100.00(0.00)
$h^{(1)}$	$H_0^{(7)}$ and $\Sigma = \mathbf{I}$					
0	5.83(0.96)	5.67(0.94)	5.83(0.96)	5.83(0.96)	5.17(0.90)	5.50(0.93)
0.4	15.33(1.47)	15.00(1.46)	15.00(1.46)	13.67(1.40)	13.17(1.38)	13.50(1.40)
0.8	36.00(1.96)	34.67(1.94)	35.50(1.95)	37.00(1.97)	35.67(1.96)	36.67(1.97)
1.6	82.67(1.55)	82.00(1.57)	82.50(1.55)	82.50(1.55)	81.17(1.60)	82.67(1.55)
	$H_0^{(7)}$ and $\Sigma = \{0.5^{ i-j }\}$					
0	6.00(0.97)	5.33(0.92)	5.50(0.93)	6.00(0.97)	5.83(0.96)	6.00(0.97)
0.4	32.83(1.92)	31.50(1.90)	32.33(1.91)	33.50(1.93)	32.67(1.91)	32.83(1.92)
0.8	78.83(1.67)	78.00(1.69)	78.83(1.68)	77.33(1.71)	76.83(1.72)	77.50(1.70)
1.6	99.50(0.29)	99.67(0.24)	99.67(0.24)	99.50(0.29)	99.83(0.17)	99.83(0.17)
$h^{(1)}$	$H_0^{(8)}$ and $\Sigma = \mathbf{I}$					
0	7.83(1.10)	7.33(1.06)	7.83(1.10)	5.83(0.96)	5.33(0.92)	5.67(0.94)
0.4	12.33(1.34)	12.00(1.33)	12.33(1.34)	12.00(1.33)	11.50(1.30)	11.67(1.31)
0.8	27.50(1.82)	26.33(1.80)	27.17(1.82)	26.00(1.79)	25.33(1.78)	26.00(1.79)
1.6	64.50(1.95)	61.00(1.99)	63.17(1.97)	66.50(1.93)	63.50(1.97)	65.67(1.94)
	$H_0^{(8)}$ and $\Sigma = \{0.5^{ i-j }\}$					
0	5.33(0.92)	4.83(0.88)	5.00(0.89)	5.50(0.93)	5.50(0.93)	5.50(0.93)
0.4	23.00(1.72)	22.00(1.70)	22.50(1.70)	25.00(1.77)	24.00(1.74)	24.33(1.75)
0.8	66.00(1.93)	63.50(1.97)	65.17(1.95)	64.17(1.97)	62.83(1.96)	63.50(1.97)
1.6	98.00(0.57)	98.17(0.55)	98.33(0.52)	97.67(0.62)	98.00(0.57)	98.00(0.57)

TABLE S7

Rejection probabilities (%) of the likelihood ratio (denoted by T_L), Wald (denoted by T_W) and score statistic (denoted by T_S) for testing $H_0^{(1)}$, with standard errors in parenthesis (%).

	$p = 50$			$p = 200$		
	T_L	T_W	T_S	T_L	T_W	T_S
$h^{(1)}$	$H_0^{(1)}$ and $\Sigma = \mathbf{I}$					
0	4.50(0.85)	4.67(0.86)	4.67(0.86)	4.17(0.82)	4.17(0.82)	4.17(0.82)
0.03	8.17(1.12)	8.17(1.12)	8.17(1.12)	10.83(1.27)	10.83(1.27)	10.83(1.27)
0.06	29.00(1.85)	29.50(1.86)	29.17(1.86)	28.50(1.84)	28.50(1.84)	28.50(1.84)
0.12	76.33(1.74)	76.67(1.73)	76.50(1.73)	73.50(1.80)	73.33(1.81)	73.50(1.80)
	$H_0^{(1)}$ and $\Sigma = \{0.5^{ i-j }\}$					
0	4.83(0.88)	4.83(0.88)	4.83(0.88)	7.17(1.05)	7.17(1.05)	7.17(1.05)
0.03	9.33(1.19)	9.50(1.20)	9.33(1.19)	11.00(1.28)	11.00(1.28)	11.00(1.28)
0.06	30.67(1.88)	30.67(1.88)	30.83(1.89)	29.50(1.86)	29.50(1.86)	29.50(1.86)
0.12	79.67(1.64)	79.67(1.64)	79.67(1.64)	77.17(1.71)	77.17(1.71)	77.17(1.71)
$h^{(2)}$	$H_0^{(2)}$ and $\Sigma = \mathbf{I}$					
0	5.33(0.92)	5.17(0.9)	5.50(0.93)	4.50(0.85)	4.50(0.85)	5.00(0.89)
0.03	19.17(1.61)	19.00(1.60)	19.67(1.62)	20.33(1.64)	20.01(1.66)	21.67(1.68)
0.06	57.33(2.02)	57.67(2.02)	57.67(2.02)	61.33(1.99)	61.67(1.98)	62.00(1.98)
0.12	99.33(0.33)	99.33(0.33)	99.33(0.33)	99.00(0.41)	99.17(0.37)	99.33(0.33)
	$H_0^{(2)}$ and $\Sigma = \{0.5^{ i-j }\}$					
0	3.83(0.78)	3.83(0.78)	3.83(0.78)	6.17(0.98)	6.17(0.98)	6.67(1.02)
0.03	16.00(1.50)	16.50(1.52)	16.33(1.51)	18.33(1.58)	18.00(1.57)	18.33(1.58)
0.06	54.50(2.03)	54.33(2.03)	54.33(2.03)	55.50(2.03)	55.00(2.03)	56.33(2.02)
0.12	99.00(0.41)	99.17(0.37)	99.17(0.37)	99.33(0.33)	99.33(0.33)	99.33(0.33)

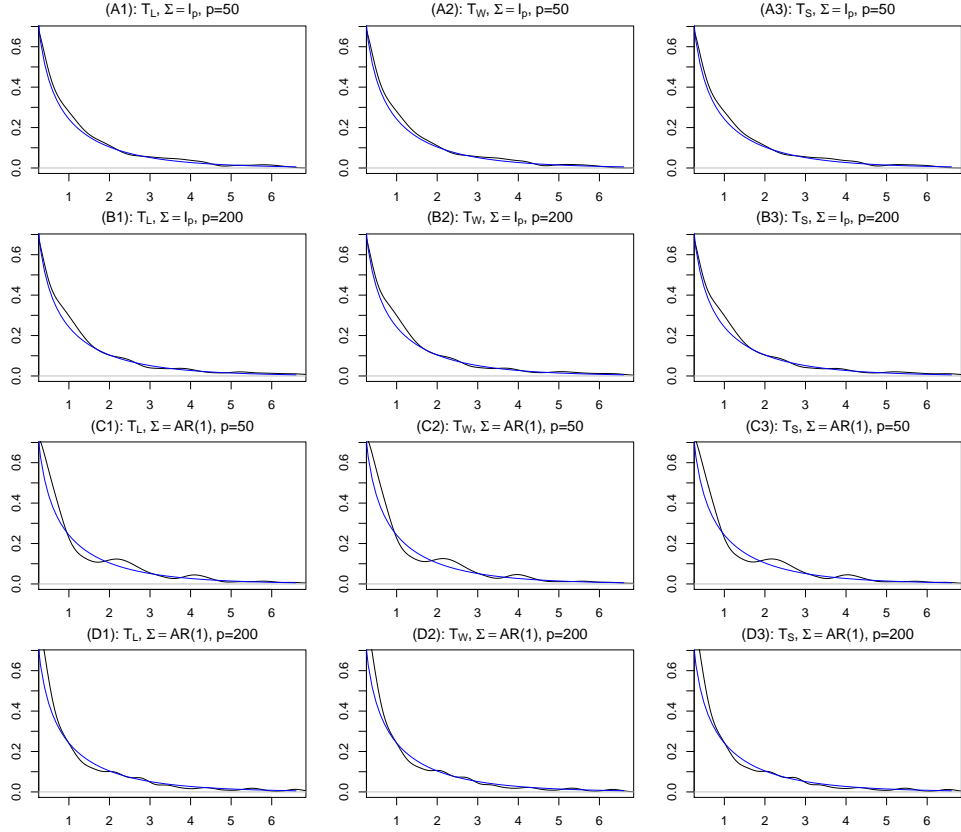


FIG S4. Kernel density plot of three test statistics under $H_0^{(1)}$ with different combinations of p and the covariance matrices. T_L , T_W and T_S from left to right. $p = 50, \Sigma = I_p$ and $p = 200, \Sigma = I_p$ and $p = 50, \Sigma = \{0.5^{|i-j|}\}_{i,j}$ and $p = 200, \Sigma = \{0.5^{|i-j|}\}_{i,j}$ from top to bottom. The black line plot the density function of a χ^2 distribution with 1 degree of freedom.

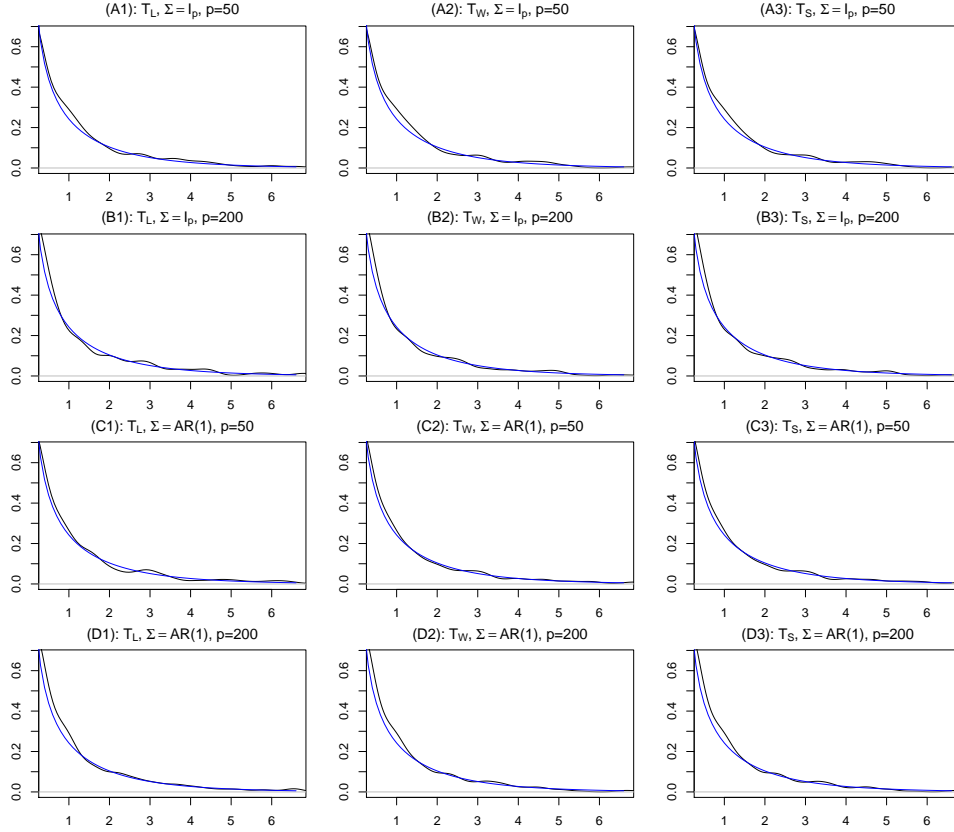


FIG S5. Kernel density plot of three test statistics under $H_0^{(2)}$ with different combinations of p and the covariance matrices. T_L , T_W and T_S from left to right. $p = 50, \Sigma = I_p$ and $p = 200, \Sigma = I_p$ and $p = 50, \Sigma = \{0.5^{|i-j|}\}_{i,j}$ and $p = 200, \Sigma = \{0.5^{|i-j|}\}_{i,j}$ from top to bottom. The black line plot the density function of a χ^2 distribution with 1 degree of freedom.

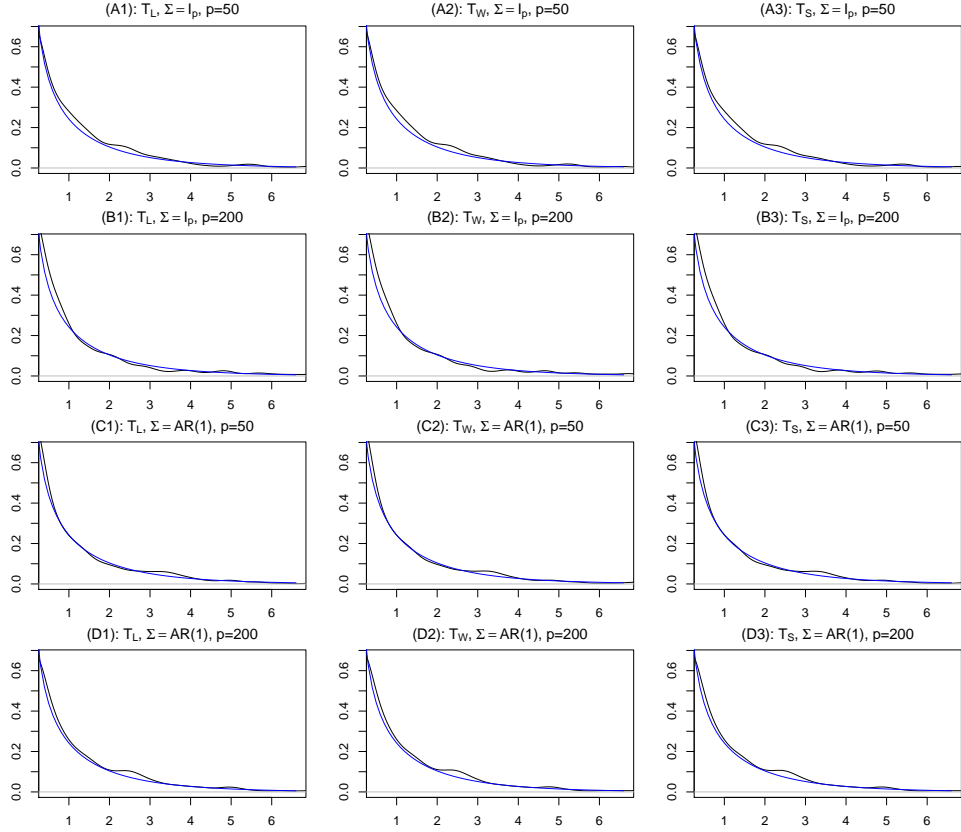


FIG S6. Kernel density plot of three test statistics under $H_0^{(2)}$ with different combinations of p and the covariance matrices. T_L , T_W and T_S from left to right. $p = 50, \Sigma = I_p$ and $p = 200, \Sigma = I_p$ and $p = 50, \Sigma = \{0.5^{|i-j|}\}_{i,j}$ and $p = 200, \Sigma = \{0.5^{|i-j|}\}_{i,j}$ from top to bottom. The black line plot the density function of a χ^2 distribution with 1 degree of freedom.

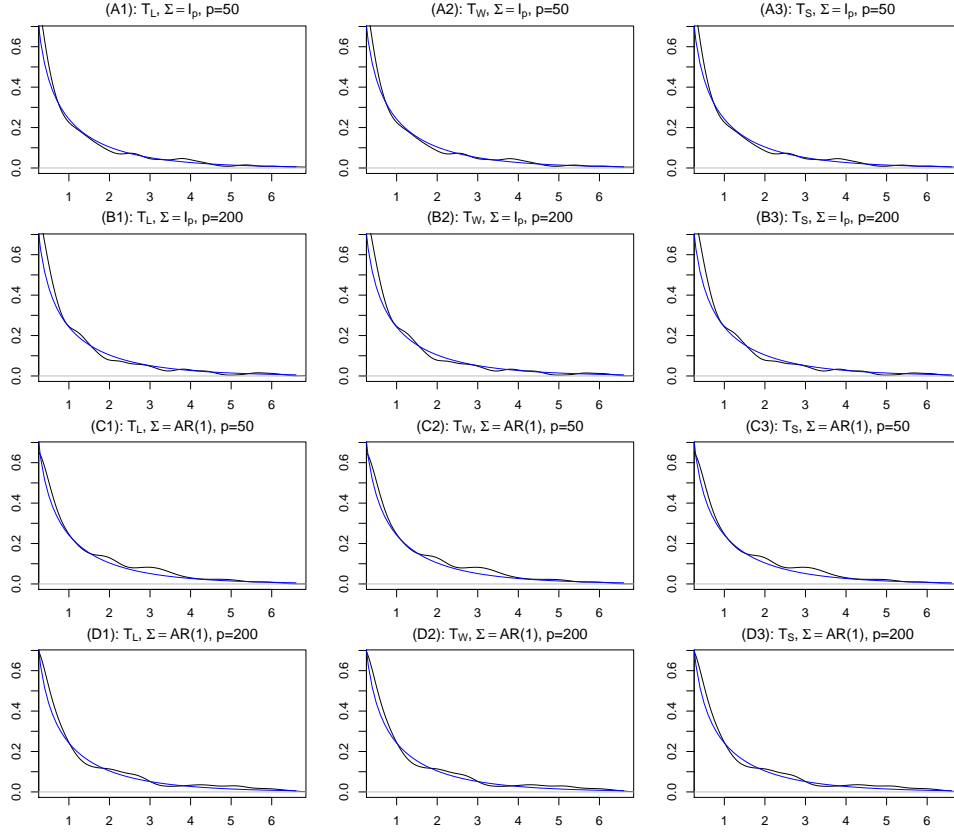


FIG S7. Histograms of three test statistics under $H_0^{(1)}$ with different combinations of p and the covariance matrices. T_L , T_W and T_S from left to right. $p = 50, \Sigma = I_p$ and $p = 200, \Sigma = I_p$ and $p = 50, \Sigma = \{0.5^{|i-j|}\}_{i,j}$ and $p = 200, \Sigma = \Sigma = \{0.5^{|i-j|}\}_{i,j}$ from top to bottom. The black line plot the density function of a χ^2 distribution with 1 degree of freedom.

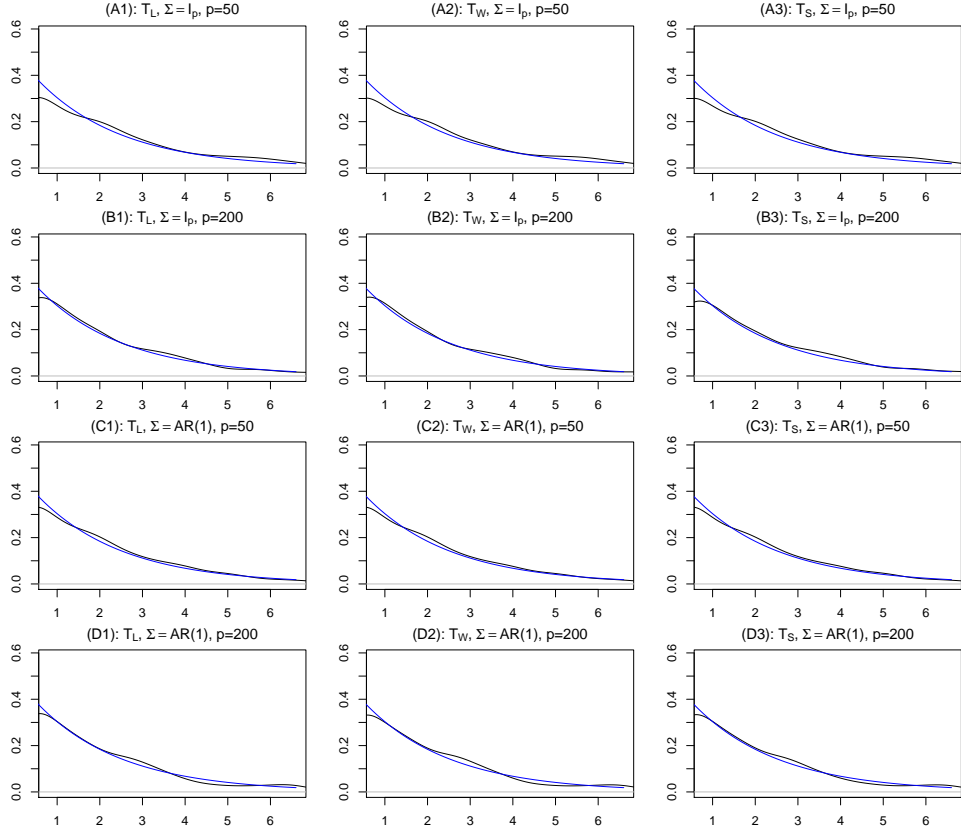


FIG S8. Histograms of three test statistics under $H_0^{(2)}$ with different combinations of p and the covariance matrices. T_L , T_W and T_S from left to right. $p = 50, \Sigma = I_p$ and $p = 200, \Sigma = I_p$ and $p = 50, \Sigma = \{0.5^{|i-j|}\}_{i,j}$ and $p = 200, \Sigma = \Sigma = \{0.5^{|i-j|}\}_{i,j}$ from top to bottom. The black line plot the density function of a χ^2 distribution with 2 degrees of freedom.

REFERENCES

- BENTKUS, V. (2004). A Lyapunov type bound in \mathbf{R}^d . *Teor. Veroyatn. Primen.* **49** 400–410.
- FAN, J. and LV, J. (2011). Nonconcave penalized likelihood with NP-dimensionality. *IEEE Trans. Inform. Theory* **57** 5467–5484.
- FAN, J. and SONG, R. (2010). Sure independence screening in generalized linear models with NP-dimensionality. *Ann. Statist.* **38** 3567–3604.
- LV, J. and FAN, Y. (2009). A Unified Approach to Model Selection and Sparse Recovery Using Regularized Least Squares. *Ann. Statist.* **37** 3498–3528.
- NING, Y. and LIU, H. (2017). A general theory of hypothesis tests and confidence regions for sparse high dimensional models. *Ann. Statist.* **45** 158–195.
- PRICE, A. L., PATTERSON, N. J., PLENGE, R. M., WEINBLATT, M. E., SHADICK, N. A. and REICH, D. (2006). Principal components analysis corrects for stratification in genome-wide association studies. *Nature genetics* **38** 904–909.
- SHI, C., FAN, A., SONG, R. and LU, W. (2017). High-dimensional A-learning for optimal dynamic treatment regimes. *Ann. Statist.* **Accepted**.
- VAN DE GEER, S., BÜHLMANN, P., RITOV, Y. and DEZEURE, R. (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *Ann. Statist.* **42** 1166–1202.

CHENGCHUN SHI AND RUI SONG
 DEPARTMENT OF STATISTICS
 NORTH CAROLINA STATE UNIVERSITY
 RALEIGH, NC 27695-8203
 USA
 E-MAIL: cshi4@ncsu.edu
 E-MAIL: rsong@ncsu.edu

ZHAO CHEN AND RUNZE LI
 DEPARTMENT OF STATISTICS,
 AND THE METHODOLOGY CENTER
 THE PENNSYLVANIA STATE UNIVERSITY,
 UNIVERSITY PARK, PA 16802-2111
 USA
 E-MAIL: zuc4@psu.edu
 E-MAIL: rzli@psu.edu