

**SUPPLEMENTARY MATERIAL TO “INTERACTION  
PURSUIT IN HIGH-DIMENSIONAL MULTI-RESPONSE  
REGRESSION VIA DISTANCE CORRELATION”**

BY YINFEI KONG<sup>¶</sup>, DAOJI LI<sup>||</sup>, YINGYING FAN<sup>\*\*</sup> AND JINCHI LV<sup>\*\*</sup>

*California State University at Fullerton<sup>¶</sup>, University of Central Florida<sup>||</sup>  
and University of Southern California<sup>\*\*</sup>*

This Supplementary Material contains some intermediate steps of the proof of Theorem 1 and additional numerical studies and technical details, as well as the details about the post-screening interaction selection.

APPENDIX B: POST-SCREENING INTERACTION SELECTION

The screening step of IPDC can reduce the problem of interaction identification from a huge scale to a moderate one as shown in Section 2. In particular, the reduced interaction model can be of dimensionality smaller than the sample size. After the screening step, IPDC further selects important interactions and main effects. Thanks to the much reduced scale, the selection step can be conducted in a computationally efficient fashion by exploiting regularization methods for the multi-response regression. Various regularization methods have been developed for multi-response linear models. See, for example, [3], [7], [5], [6], [25], and references therein. Those methods were usually investigated for the scenario of no interactions. For the selection step of IPDC, we aim at interaction model recovery by employing a two-step variable selection procedure, where we first recover the support union using the idea of group variable selection [35] and then estimate the individual supports for each column of the regression coefficient matrix via an additional refitting step of Lasso [32] applied to the recovered support union (see, e.g., [14] for connections and differences among regularization methods).

To simplify the presentation, hereafter we assume that the response vector  $\mathbf{y}$  is centered with mean zero and all interactions  $X_k X_\ell$  are also centered to have mean zero with a slight abuse of notation, which eliminates the intercept vector  $\boldsymbol{\alpha}$  in model (1). Thus given an i.i.d. sample  $(\mathbf{y}_i, \mathbf{x}_i)_{i=1}^n$ , the multi-response interaction model (1) can be rewritten in the matrix form

$$(A.1) \quad \mathbf{Y} = \tilde{\mathbf{X}}\mathbf{B} + \mathbf{W},$$

where  $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)^T \in \mathbb{R}^{n \times q}$  is the response matrix,  $\tilde{\mathbf{X}} = (\mathbf{X}, \mathbf{Z}) \in \mathbb{R}^{n \times \tilde{p}}$  with  $\tilde{p} = p(p+1)/2$  is the full augmented design matrix with  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T \in \mathbb{R}^{n \times p}$  the main effect matrix and  $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)^T \in \mathbb{R}^{n \times [p(p-1)/2]}$  the interaction matrix,  $\mathbf{z}_i$ 's are defined similarly to  $\mathbf{z}$ ,  $\mathbf{B} = (\mathbf{B}_x^T, \mathbf{B}_z^T)^T \in \mathbb{R}^{\tilde{p} \times q}$  is the regression coefficient matrix with  $\mathbf{B}_x \in \mathbb{R}^{p \times q}$  and  $\mathbf{B}_z \in \mathbb{R}^{[p(p-1)/2] \times q}$ , and  $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_n)^T \in \mathbb{R}^{n \times q}$  is the error matrix.

**B.1. Interaction and main effect selection.** Let  $S$  be the row support of the true regression coefficient matrix  $\mathbf{B}^*$  in model (A.1), which corresponds to the index set of nonzero rows of  $\mathbf{B}^*$ ; that is, if  $k \in S$ , then the  $k$ th row of  $\mathbf{B}^*$  has at least one nonzero component. Denote by  $\mathcal{J} = \{j_1, \dots, j_{d_1}\} \subset \{1, \dots, p\}$  and  $\mathcal{K} = \{k_1, \dots, k_{d_2}\} \subset \{1, \dots, p\}$  the index sets of retained main effects and interaction variables after the screening step of IPDC, respectively. Then the reduced design matrix is  $(\tilde{\mathbf{x}}_{j_1}, \dots, \tilde{\mathbf{x}}_{j_{d_1}}, \tilde{\mathbf{x}}_{k_1} \circ \tilde{\mathbf{x}}_{k_2}, \dots, \tilde{\mathbf{x}}_{k_{d_2-1}} \circ \tilde{\mathbf{x}}_{k_{d_2}}) \in \mathbb{R}^{n \times d}$  with  $d = d_1 + d_2(d_2 - 1)/2$ , where  $\tilde{\mathbf{x}}_\ell$  is the  $\ell$ th column of  $\mathbf{X}$ . Let  $\tilde{S} \subset \{1, \dots, \tilde{p}\}$  be the index set given by the columns of such a reduced matrix in the full matrix  $\tilde{\mathbf{X}}$ . As guaranteed by Theorem 1, the true row support  $S$  can be contained in the reduced set  $\tilde{S}$  by IPDC with high probability that converges to one at a fast rate as sample size  $n$  increases.

Observe that the true row support  $S$  is the union of individual supports of the columns of the true regression coefficient matrix  $\mathbf{B}^*$  corresponding to the  $q$  responses. Given any set  $J \subset \{1, \dots, \tilde{p}\}$ , denote by  $\mathbf{B}_J$  a submatrix of  $\mathbf{B}$  formed by the rows indexed by  $J$ . For the support union recovery, we exploit the multivariate group Lasso given by the following regularization problem

$$(A.2) \quad \min_{\mathbf{B}_{\tilde{S}^c} = \mathbf{0}} \left\{ \frac{1}{2nq} \|\mathbf{Y} - \tilde{\mathbf{X}}\mathbf{B}\|_F^2 + \lambda \|\mathbf{B}\|_{2,1} \right\},$$

where  $\tilde{S}^c$  is the complement of the set  $\tilde{S}$ ,  $\|\cdot\|_F$  denotes the Frobenius norm of a matrix,  $\lambda \geq 0$  is a regularization parameter, and  $\|\cdot\|_{2,1}$  stands for the matrix rowwise (2,1)-norm defined as  $\|\mathbf{M}\|_{2,1} = \sum_i (\sum_j m_{ij}^2)^{1/2}$  for any matrix  $\mathbf{M} = (m_{ij})$ . Note that  $\tilde{S}$  and  $\widehat{\mathcal{M}} \cup \widehat{\mathcal{I}}$  share the same cardinality. We should remark that as ensured by Theorem 2, the computational cost of solving the optimization problem (A.2) can be substantially reduced compared to that of solving the same optimization problem without the screening step, that is, with  $\tilde{S} = \{1, \dots, \tilde{p}\}$ .

The multivariate group Lasso has been widely used in the multi-response linear regression models typically without interaction terms. For example,

[24] and [25] established the oracle inequalities for the case when the design matrix is deterministic and the error matrix has i.i.d. Gaussian entries. [28] investigated the model selection consistency in terms of support union recovery of the multivariate group Lasso under the assumptions that the design matrix is drawn with i.i.d. Gaussian row vectors and all the entries of the error matrix are i.i.d. Gaussian. We will relax such Gaussianity assumptions and justify that this group variable selection procedure continues to perform well in the presence of interactions.

Once the row support of the true regression coefficient matrix is recovered, it is straightforward to recover the individual supports of each column of the regression coefficient matrix by an additional refitting step of applying the ordinary Lasso to the recovered support union. Since the sampling properties of Lasso have been extensively studied and are now well understood in the literature, we will provide only theoretical analysis of the multivariate group Lasso problem (A.2).

**B.2. Support union recovery and oracle inequalities.** To facilitate our technical analysis for the selection step of IPDC, we impose a few additional regularity conditions.

CONDITION 4. *The covariate vector  $\mathbf{x}$  has a sub-Gaussian distribution and  $s = |S| = O(n^\xi)$  for some constant  $0 \leq \xi < 1/4$ .*

CONDITION 5 (RE( $s$ ) assumption). *There exists some positive constant  $\kappa$  such that*

$$\kappa(s) = \min_{|J| \leq s, \mathbf{\Delta} \in \mathbb{R}^{\tilde{p} \times q} \setminus \{\mathbf{0}\}, \|\mathbf{\Delta}_{J^c}\|_{2,1} \leq 3\|\mathbf{\Delta}_J\|_{2,1}} \frac{\|\mathbf{\Sigma}^{1/2} \mathbf{\Delta}\|_F}{\|\mathbf{\Delta}_J\|_F} \geq \kappa,$$

where  $\mathbf{\Sigma}$  is the covariance matrix of  $\tilde{\mathbf{x}} = (\mathbf{x}^T, \mathbf{z}^T)^T$ .

CONDITION 6. *The error vector  $\mathbf{w}$  has a sub-exponential distribution.*

The first part of Condition 4 is a mild assumption on the distribution of the covariates. It can be satisfied by many light-tailed distributions such as Gaussian distributions and distributions with bounded support. The second part of Condition 4 puts a row sparsity constraint on the true regression coefficient matrix. In particular, the requirement of  $\xi < 1/4$  reflects the difficulty of interaction selection in high dimensions.

Condition 5 is a natural extension of the restricted eigenvalue (RE) assumption introduced in [1] since here we use the rowwise  $(2, 1)$ -norm in

place of the  $L_1$ -norm. The RE assumption has been commonly used to establish the oracle inequalities for the Lasso and Dantzig selector [4]. For simplicity we still refer to Condition 5 as the RE( $s$ ) assumption. This condition is also similar to Condition 3.1 in [24] and Condition 4.1 in [25], who considered the scenario of deterministic design matrix and no interactions. Condition 6 assumes the sub-exponential distribution for the error vector, which is key to establishing the deviation probability bound for  $\|\tilde{\mathbf{X}}^T \mathbf{W}\|_{2,\infty}$ . Here  $\|\cdot\|_{2,\infty}$  denotes the matrix rowwise  $(2, \infty)$ -norm defined as  $\|\mathbf{M}\|_{2,\infty} = \max_i (\sum_j m_{ij}^2)^{1/2}$  for any matrix  $\mathbf{M} = (m_{ij})$ . Hereafter  $p$  involved in the regularization parameter  $\lambda$  and probability bounds is understood implicitly as  $\max\{n, p\}$ .

**THEOREM 3.** *Assume that all the conditions of Theorem 1 and Conditions 4–6 hold,  $q \leq p$ ,  $\log p = o(n^\eta)$  with  $\eta = \min\{\eta_0, 1/2 - 2\xi\}$ , and set  $\lambda = c_3 \sqrt{(\log p)/(nq)}$  with  $c_3 > 0$  some constant. Then with probability at least  $1 - O\{\exp(-Cn^{\eta_0/2})\} - O(p^{-c_4})$  for some constants  $C, c_4 > 0$ , the minimizer  $\hat{\mathbf{B}}$  of (A.2) satisfies*

$$(A.3) \quad (nq)^{-1/2} \|\tilde{\mathbf{X}}(\hat{\mathbf{B}} - \mathbf{B}^*)\|_F \leq \frac{8c_3}{\kappa} \sqrt{s(\log p)/n},$$

$$(A.4) \quad \frac{1}{\sqrt{q}} \|\hat{\mathbf{B}} - \mathbf{B}^*\|_{2,1} \leq \frac{64c_3}{\kappa^2} s \sqrt{(\log p)/n}.$$

*If in addition  $\min_{j \in S} \|\mathbf{B}_j^*\|/\sqrt{q} > 128c_3\kappa^{-2}s\sqrt{(\log p)/n}$ , then with the same probability the row support of  $\tilde{\mathbf{B}}$  is identical to  $S$ , where the matrix  $\tilde{\mathbf{B}}$  is obtained by thresholding the  $j$ th row of  $\hat{\mathbf{B}}$  to zero for each  $j$  if  $\|\hat{\mathbf{B}}_j\|/\sqrt{q} \leq 64c_3\kappa^{-2}s\sqrt{(\log p)/n}$ . Moreover, if the RE( $s$ ) assumption in Condition 5 is replaced by RE( $2s$ ), then it holds with the same probability that*

$$(A.5) \quad \frac{1}{\sqrt{q}} \|\hat{\mathbf{B}} - \mathbf{B}^*\|_F \leq \frac{16\sqrt{10}c_3}{\kappa^2(2s)} \sqrt{s(\log p)/n}.$$

Theorem 3 establishes the model selection consistency of the IPDC followed by hard thresholding in terms of support union recovery. It also extends the oracle inequalities in Theorem 3.3 of [24] and Corollary 4.1 of [25] in three important aspects: the inclusion of interaction terms, the analysis of large random design matrix, and the relaxed distributional assumption. Such extensions make the technical analyses more involved and challenging. We should remark that the same results as in Theorem 3 hold with probability at least  $1 - O(p^{-c_4})$  for the regularized estimator with  $d = \tilde{p}$ , that is, without the screening step. It is also worth mentioning that the value

TABLE 5

*Proportions of important main effects, important interaction, and all of them retained by different screening methods.*

Method	$X_{12}$	$X_{22}$	$X_1X_2$	All
DCSIS2	1.00	1.00	0.06	0.06
DCSIS-square	0.00	0.81	1.00	0.00
IPDC	1.00	1.00	1.00	1.00

of the regularization parameter  $\lambda$  in our Theorem 3 is slightly larger than those used in [24] and [25], due to the more general model setting considered in our paper. In fact, such larger value of  $\lambda$  is needed to suppress the additional noise caused by the presence of interactions and the heavier-tailed distribution of model errors.

## APPENDIX C: ADDITIONAL NUMERICAL STUDIES

**C.1. Comparison of IPDC with individual components.** Recall that the new interaction screening approach of IPDC treats the screening for interactions and the screening for main effects as two separate components. Since the distance correlation can capture nonlinear dependency between variables, a natural question is whether either of these two components might suffice for the purpose of the joint screening for both interactions and main effects. To ease the presentation, the component for interaction screening is referred to as DCSIS-square, and the component for main effect screening is called DCSIS2 as described in Section 3.1. Thus it is of interest to compare IPDC with both DCSIS2 and DCSIS-square. To this end, we revisit the setting 3 of Model 4 investigated in Section 3.1; see Table 1 for the screening performance of DCSIS2 and IPDC.

Table 5 reports the comparison results of these three methods. We see that DCSIS2, which is designed specifically for main effect screening, fails to retain the important interaction  $X_1X_2$ , and DCSIS-square, which is designed specifically for interaction screening, fails to retain the important main effect  $X_{12}$ . In contrast, the IPDC combines the strengths of its two individual components in screening for both interactions and main effects. Such an observation is in line with a key message spelled out in the Introduction, that is, a separate screening step for interactions can significantly enhance the screening performance if one aims at finding important interactions.

## C.2. Performance of interaction and main effect selection.

C.2.1. *Selection in single-response models.* After the screening step, we further investigate the performance of selection for the interactions and main effects in the reduced feature space. The selection step in the single-response examples (Models 1–4) is implemented by the Lasso. Thus we refer to each two-stage interaction screening and selection procedure by the SIS2-Lasso, DCSIS2-Lasso, SIRI-Lasso, IP-Lasso, and IPDC-Lasso, respectively. The oracle procedure, which assumes that the true underlying sparse interaction model is known in advance, is used as a benchmark for comparison. In particular, in Model 4 the indicator covariate  $\mathbb{I}(X_{12} \geq 0)$  instead of the linear predictor  $X_{12}$  is used in the oracle procedure.

Three performance measures, the prediction error (PE), the number of false positives (FP), and the number of false negatives (FN), are employed to assess the variable selection performance of each method in the single-response examples. The PE is defined as  $E(Y - \hat{Y})^2$  with  $\hat{Y}$  the predicted response. We generate an independent test sample of size 10,000 to calculate the PE. The FP is defined as the total number of unimportant interactions and main effects included in the final model, while the FN is defined as the total number of important interactions and main effects missed by the final model.

Table 6 presents the means and standard errors of different performance measures. Since the screening results by all methods in Model 1 under three different settings are almost identical, one can expect that the corresponding selection results should be very similar, which is indeed the case. Thus we omit the selection results for Model 1 to save space. For Model 2, the performance of all methods is very similar across three settings. As a result, we only present the selection results under setting 2 of this model in Table 7 as a representative. The complete results are available upon request. Based on Tables 6 and 7, the following observations can be made.

- In Model 2, we see that IPDC-Lasso performs the best and closest to the oracle procedure across all measures, as shown in Table 7. For Model 3 under all settings, our method IPDC-Lasso has far lower mean prediction error than all other methods except for the oracle, according to Table 6. The advantage of IPDC-Lasso over other methods is also evident in Model 4.
- We remark that the gap between the prediction errors of the IPDC-Lasso and oracle in Model 4 is mainly because as mentioned before, the latter exploits the indicator covariate  $\mathbb{I}(X_{12} \geq 0)$  whereas such prior information is unavailable to all other procedures. Even in this scenario of model misspecification, our method still performs well in identifying important interactions and main effects.

TABLE 6  
*Means and standard errors (in parentheses) of different selection performance measures for Models 3 and 4 over 100 replications.*

Method	Model 3			Model 4		
	PE	FP	FN	PE	FP	FN
Setting 1: $(p, \rho) = (2000, 0.5)$						
SIS2-Lasso	33.79 (0.87)	39.25 (3.79)	1.88 (0.04)	15.00 (0.32)	22.05 (2.50)	1.23 (0.06)
DCSIS2-Lasso	3.94 (0.41)	0.51 (0.17)	0.32 (0.05)	3.36 (0.38)	4.16 (0.72)	0.11 (0.04)
SIRI-Lasso	3.54 (0.40)	0.54 (0.36)	0.27 (0.05)	4.34 (0.49)	1.63 (0.40)	0.31 (0.07)
IP-Lasso	2.08 (0.28)	0.46 (0.11)	0.10 (0.03)	2.38 (0.05)	4.24 (0.64)	0.07 (0.03)
IPDC-Lasso	1.27 (0.10)	0.63 (0.20)	0.01 (0.01)	2.27 (0.02)	3.32 (0.47)	0.01 (0.01)
Oracle	1.017 (0.002)	0 (0)	0 (0)	1.022 (0.002)	0 (0)	0 (0)
Setting 2: $(p, \rho) = (5000, 0.5)$						
SIS2-Lasso	36.31 (0.51)	61.78 (2.85)	1.97 (0.02)	15.27 (0.24)	39.29 (2.33)	1.15 (0.04)
DCSIS2-Lasso	5.81 (0.44)	1.20 (0.55)	0.54 (0.05)	4.17 (0.48)	3.45 (0.53)	0.20 (0.05)
SIRI-Lasso	4.48 (0.45)	0.42 (0.16)	0.37 (0.05)	4.70 (0.52)	2.24 (0.44)	0.37 (0.07)
IP-Lasso	2.52 (0.33)	1.83 (0.61)	0.15 (0.04)	2.51 (0.06)	6.75 (1.24)	0.14 (0.04)
IPDC-Lasso	1.38 (0.16)	0.91 (0.31)	0.02 (0.01)	2.30 (0.02)	4.39 (0.63)	0.01 (0.01)
Oracle	1.009 (0.002)	0 (0)	0 (0)	1.014 (0.002)	0 (0)	0 (0)
Setting 3: $(p, \rho) = (2000, 0.1)$						
SIS2-Lasso	21.96 (0.19)	22.84 (3.18)	1.98 (0.01)	13.15 (0.10)	15.68 (2.04)	1.24 (0.05)
DCSIS2-Lasso	18.85 (0.47)	9.32 (2.47)	1.70 (0.05)	12.58 (0.28)	8.51 (1.73)	1.15 (0.05)
SIRI-Lasso	14.45 (0.72)	0.40 (0.17)	1.28 (0.07)	11.55 (0.45)	1.35 (0.43)	1.28 (0.07)
IP-Lasso	6.23 (0.63)	4.20 (1.37)	0.46 (0.06)	2.54 (0.17)	6.28 (1.54)	0.05 (0.02)
IPDC-Lasso	3.08 (0.44)	0.99 (0.21)	0.17 (0.04)	2.26 (0.01)	4.00 (0.79)	0.00 (0.00)
Oracle	1.017 (0.002)	0 (0)	0 (0)	1.022 (0.002)	0 (0)	0 (0)

TABLE 7  
*Means and standard errors (in parentheses) of different selection performance measures for setting 2 of Model 2 over 100 replications.*

Method	PE	FP	FN
SIS2-Lasso	25.57 (1.61)	30.20 (3.31)	1.62 (0.10)
DCSIS2-Lasso	3.20 (0.40)	1.85 (0.44)	0.21 (0.04)
SIRI-Lasso	3.03 (0.38)	1.30 (0.23)	0.20 (0.04)
IP-Lasso	4.05 (0.45)	4.79 (1.06)	0.33 (0.05)
IPDC-Lasso	1.61 (0.20)	2.55 (0.49)	0.04 (0.02)
Oracle	1.014 (0.002)	0 (0)	0 (0)

TABLE 8  
*Means and standard errors (in parentheses) of different selection performance measures for Model 5 over 100 replications.*

Method	PE	FP.main	FP.int	FN.main	FN.int
SIS.max-GLasso	6.24 (0.06)	127.64 (4.72)	699.08 (27.37)	1.22 (0.23)	9.50 (0.12)
SIS.sum-GLasso	6.12 (0.08)	185.28 (3.64)	810.72 (20.19)	0.74 (0.14)	9.06 (0.17)
DCSIS-GLasso	4.29 (0.14)	157.74 (3.98)	830.19 (23.00)	0.50 (0.10)	5.63 (0.29)
IPDC-GLasso	2.74 (0.09)	124.77 (3.02)	878.93 (20.60)	0.04 (0.02)	2.46 (0.23)
SIS.max-GLasso-Lasso	5.11 (0.05)	11.92 (0.84)	52.40 (3.17)	3.44 (0.19)	9.60 (0.10)
SIS.sum-GLasso-Lasso	4.99 (0.07)	15.75 (0.86)	63.73 (3.39)	3.07 (0.20)	9.35 (0.16)
DCSIS-GLasso-Lasso	3.40 (0.12)	11.87 (0.70)	62.98 (2.94)	1.41 (0.17)	6.36 (0.28)
IPDC-GLasso-Lasso	2.08 (0.09)	7.42 (0.36)	58.95 (2.32)	0.37 (0.09)	3.28 (0.24)
Oracle	1.048 (0.001)	0 (0)	0 (0)	0 (0)	0 (0)

*C.2.2. Selection in multi-response model.* As mentioned before, interaction and main effect selection in the multi-response model setting (Model 5) is conducted through a two-step procedure in the reduced feature space obtained by screening. Such a method first selects rows of the regression coefficient matrix using the group Lasso, and then applies the Lasso to each individual response for further selection of the rows. The goal of the individual Lasso is to eliminate the unimportant interactions and main effects that are included in the model recovered by the group Lasso. The resulting interaction screening and selection procedures are referred to as the SIS.max-GLasso-Lasso, SIS.sum-GLasso-Lasso, DCSIS-GLasso-Lasso, and IPDC-GLasso-Lasso, respectively. We also include for comparison the procedures that exploit only the group Lasso in the selection step, which are named by dropping the Lasso component. The oracle procedure is the ordinary least-squares estimation applied to each response separately on the corresponding true support.

The same performance measures as defined in Section C.2.1 are employed to evaluate different methods, except that the PE is now calculated as the average prediction error across all  $q = 10$  responses. To further differentiate the false positives and false negatives for the main effects and interactions, we attach “.main” and “.int” to both measures of FP and FN as shown in Table 8.

Table 8 reports the selection results for Model 5. The FP.int is relatively large for all methods since even after screening, there are still a large number of interactions left, due to the presence of multiple responses. We observe that a further step of individual Lasso implemented on the support of group lasso for each response separately can substantially reduce the FP for both interactions and main effects. Moreover, our method IPDC-GLasso-Lasso outperforms all competitors under all performance measures.



TABLE 9

Means and standard errors (in parentheses) of prediction error as well as numbers of selected main effects and interactions for each method in mice data.

Method	PE ( $\times 10^{-3}$ )	Model Size	
		Main	Interaction
SIS2-Lasso	100.14 (5.57)	0.14 (0.03)	9.56 (0.11)
DCSIS2-Lasso	100.91 (6.12)	0.12 (0.03)	9.22 (0.11)
SIRI-Lasso	247.53 (10.11)	1.07 (0.03)	0.46 (0.13)
IP-Lasso	101.34 (4.99)	5.00 (0.09)	8.55 (0.21)
IPDC-Lasso	96.55 (5.40)	3.04 (0.07)	7.31 (0.11)

**C.3. Univariate gene expression study.** We study the the inbred mouse microarray gene expression data set in Lan et al. [21]. There are 60 mouse arrays, with 31 from female mice and 29 from male mice. The response variable is the gene expression level of stearoyl-CoA desaturase-1 (Scd1), a gene involved in fat storage. Specifically, Scd1 controls lipid metabolism and insulin sensitivity. The covariates are gene expression levels for 22,690 of the mice’s other genes. Therefore, the sample size  $n = 60$ , the number of covariates  $p = 22,690$ , and the number of responses  $q = 1$ . All variables involved in this study are continuous.

This data set is publicly available on the Gene Expression Omnibus website (<http://www.ncbi.nlm.nih.gov/geo>; accession number GSE3330), and has been studied in Hao and Zhang [16] and Narisetty and He [27]. Following Narisetty and He [27], we randomly split the data into training and test sets of sizes 55 and 5, respectively. Furthermore, to ameliorate the numerical instability caused by the relatively small sample size we perform 200 random splits and calculate the mean prediction errors and the corresponding standard errors to better evaluate the performance of various methods. We compare the IPDC with the SIS2, DCSIS2, SIRI, and IP. Detailed descriptions of all these methods can be found in Sections 3 and C.2.

The final selection results on the prediction error and selected model size are summarized in Table 9. We see that IPDC-Lasso performs noticeably better than its competitors. Further, paired  $t$ -tests of prediction errors on the 200 splits of IPDC-Lasso against SIS2-Lasso, DCSIS2-Lasso, SIRI-Lasso, and IP-Lasso lead to  $p$ -values  $3.40 \times 10^{-2}$ ,  $1.94 \times 10^{-2}$ ,  $1.90 \times 10^{-36}$ , and  $2.64 \times 10^{-2}$ , respectively. These test results demonstrate the significantly improved performance of IPDC over other methods at the 5% level.

#### APPENDIX D: ADDITIONAL PROOFS OF MAIN RESULTS

**D.1. Step 1.3 of Part 1 in the proof of Theorem 1.** We now consider the term  $\widehat{T}_{k1,3}^* - T_{k1,3}$ . Applying the Cauchy-Schwarz inequality twice leads to

$$(A.6) \quad T_{k1,3} \leq \{E[\psi^2(\mathbf{y}_1^*, \mathbf{y}_2^*)] E[\phi^2(X_{1k}^*, X_{2k}^*) \mathbb{I}\{\phi(X_{1k}^*, X_{2k}^*) > M_1\}]\}^{1/2} \\ \leq E^{1/2}[\psi^2(\mathbf{y}_1^*, \mathbf{y}_2^*)] \{E[\phi^4(X_{1k}^*, X_{2k}^*)] P\{\phi(X_{1k}^*, X_{2k}^*) > M_1\}\}^{1/4}.$$

In view of (17), we have

$$(A.7) \quad E[\phi^4(X_{1k}^*, X_{2k}^*)] \leq E[(X_{1k}^2 + X_{2k}^2)^4] \leq E\{[2(X_{1k}^4 + X_{2k}^4)]^2\} \\ \leq E[8(X_{1k}^8 + X_{2k}^8)] = 16E(X_{1k}^8)$$

and by Bonferroni's inequality,

$$(A.8) \quad P\{\phi(X_{1k}^*, X_{2k}^*) > M_1\} \leq P(X_{1k}^2 + X_{2k}^2 > M_1) \leq P(X_{1k}^2 > M_1/2) \\ + P(X_{2k}^2 > M_1/2) = 2P(X_{1k}^2 > M_1/2) \\ \leq 2 \exp(-c_0 M_1/2) E[\exp(c_0 X_{1k}^2)].$$

Combining (25) with (A.6)–(A.8) and by Condition 2, we obtain  $T_{k1,3} \leq \widetilde{C}_3 \exp(-8^{-1} c_0 M_1)$ , where  $\widetilde{C}_3$  is some positive constant. Since  $M_1 = n^{\xi_1}$ , it holds that for any positive constant  $\widetilde{C}$ ,

$$(A.9) \quad 0 \leq T_{k1,3} \leq \widetilde{C}_3 \exp(-8^{-1} c_0 n^{\xi_1}) \leq \widetilde{C} n^{-\kappa_2} / 48$$

for all  $1 \leq k \leq p$  when  $n$  is sufficiently large. This entails that

$$(A.10) \quad P(\max_{1 \leq k \leq p} |\widehat{T}_{k1,3}^* - T_{k1,3}| \geq \widetilde{C} n^{-\kappa_2} / 24) \\ \leq P(\max_{1 \leq k \leq p} |\widehat{T}_{k1,3}^*| \geq \widetilde{C} n^{-\kappa_2} / 48)$$

for all  $n$  sufficiently large. Thus applying Markov's inequality and noting that  $\widehat{T}_{k1,3}^* \geq 0$  and  $E(\widehat{T}_{k1,3}^*) = T_{k1,3}$ , we have

$$P(|\widehat{T}_{k1,3}^*| \geq \delta/2) \leq (\delta/2)^{-1} E(|\widehat{T}_{k1,3}^*|) \leq (\delta/2)^{-1} E(\widehat{T}_{k1,3}^*) \\ = (\delta/2)^{-1} T_{k1,3}$$

for any  $\delta > 0$ . Choosing  $\delta = \widetilde{C} n^{-\kappa_2} / 24$  in the above inequality and in view of (A.9), it follows that  $P(|\widehat{T}_{k1,3}^*| \geq \widetilde{C} n^{-\kappa_2} / 48) \leq 48 \widetilde{C}^{-1} \widetilde{C}_3 n^{\kappa_2} \exp(-8^{-1} c_0 n^{\xi_1})$ . This inequality together with (A.10) and Bonferroni's inequality yields

$$(A.11) \quad P(\max_{1 \leq k \leq p} |\widehat{T}_{k1,3}^* - T_{k1,3}| \geq \widetilde{C} n^{-\kappa_2} / 24) \\ \leq 48p \widetilde{C}^{-1} \widetilde{C}_3 n^{\kappa_2} \exp(-8^{-1} c_0 n^{\xi_1}).$$

**D.2. Step 2 of Part 1 in the proof of Theorem 1.** In this step, we handle the term  $\max_{1 \leq k \leq p} |\widehat{T}_{k2} - T_{k2}|$ . Define  $T_{k2,1} = E[\phi(X_{1k}^*, X_{2k}^*)]$  and  $T_{k2,2} = E[\psi(\mathbf{y}_1^*, \mathbf{y}_2^*)]$ . Then  $T_{k2} = T_{k2,1}T_{k2,2}$ . Similarly,  $\widehat{T}_{k2}$  can be rewritten as  $\widehat{T}_{k2} = \widehat{T}_{k2,1}\widehat{T}_{k2,2}$  by letting  $\widehat{T}_{k2,1} = n^{-2} \sum_{i,j=1}^n \phi(X_{ik}^*, X_{jk}^*)$  and  $\widehat{T}_{k2,2} = n^{-2} \sum_{i,j=1}^n \psi(\mathbf{y}_i^*, \mathbf{y}_j^*)$ . An application of similar arguments as in Step 1 results in that for any positive constant  $\widetilde{C}$ , there exist some positive constants  $\widetilde{C}_1, \dots, \widetilde{C}_4$  such that

$$\begin{aligned} P(|\widehat{T}_{k2,1} - T_{k2,1}| \geq \widetilde{C}n^{-\kappa_2}/4) &\leq \widetilde{C}_1 \exp\{-\widetilde{C}_2 n^{(1-2\kappa_2)/3}\}, \\ P(|\widehat{T}_{k2,2} - T_{k2,2}| \geq \widetilde{C}n^{-\kappa_2}/4) &\leq \widetilde{C}_3 \exp\{-\widetilde{C}_4 n^{(1-2\kappa_2)/5}\}. \end{aligned}$$

In view of (17) and (18), we have  $T_{k2,1} \leq 2E(X_{1k}^2)$  and  $T_{k2,2} \leq 2E(\|\widetilde{\mathbf{y}}_1\|^2)$ . By Condition 2,  $T_{k2,1}$  and  $T_{k2,2}$  are uniformly bounded from above by some positive constant for all  $1 \leq k \leq p$ . Thus it follows from Lemma 1 that for any positive constant  $\widetilde{C}$ , there exist some positive constants  $\widetilde{C}_5$  and  $\widetilde{C}_6$  such that

$$\begin{aligned} P(|\widehat{T}_{k2} - T_{k2}| \geq \widetilde{C}n^{-\kappa_2}/4) &= P(|\widehat{T}_{k2,1}\widehat{T}_{k2,2} - T_{k2,1}T_{k2,2}| \geq \widetilde{C}n^{-\kappa_2}/4) \\ &\leq \widetilde{C}_5 \exp\{-\widetilde{C}_6 n^{(1-2\kappa_2)/5}\} \end{aligned}$$

for all  $1 \leq k \leq p$ . By Bonferroni's inequality, we obtain

$$\begin{aligned} \text{(A.12)} \quad P\left(\max_{1 \leq k \leq p} |\widehat{T}_{k2} - T_{k2}| \geq \widetilde{C}n^{-\kappa_2}/4\right) &\leq \sum_{k=1}^p P(|\widehat{T}_{k2} - T_{k2}| \geq \widetilde{C}n^{-\kappa_2}/4) \\ &\leq p\widetilde{C}_5 \exp\{-\widetilde{C}_6 n^{(1-2\kappa_2)/5}\}. \end{aligned}$$

**D.3. Step 3 of Part 1 in the proof of Theorem 1.** We now consider the term  $\widehat{T}_{k3} - T_{k3}$ . Define a  $U$ -statistic

$$\widehat{T}_{k3}^* = 6[n(n-1)(n-2)]^{-1} \sum_{i < j < l} g(X_{ik}^*, \mathbf{y}_i^*; X_{jk}^*, \mathbf{y}_j^*; X_{lk}^*, \mathbf{y}_l^*)$$

with the kernel  $g(X_{ik}^*, \mathbf{y}_i^*; X_{jk}^*, \mathbf{y}_j^*; X_{lk}^*, \mathbf{y}_l^*)$  given by

$$\begin{aligned} g(X_{ik}^*, \mathbf{y}_i^*; X_{jk}^*, \mathbf{y}_j^*; X_{lk}^*, \mathbf{y}_l^*) &= \phi(X_{ik}^*, X_{jk}^*)\psi(\mathbf{y}_i^*, \mathbf{y}_l^*) + \phi(X_{ik}^*, X_{lk}^*)\psi(\mathbf{y}_i^*, \mathbf{y}_j^*) \\ &\quad + \phi(X_{jk}^*, X_{ik}^*)\psi(\mathbf{y}_j^*, \mathbf{y}_l^*) + \phi(X_{jk}^*, X_{lk}^*)\psi(\mathbf{y}_j^*, \mathbf{y}_i^*) \\ &\quad + \phi(X_{lk}^*, X_{ik}^*)\psi(\mathbf{y}_l^*, \mathbf{y}_j^*) + \phi(X_{lk}^*, X_{jk}^*)\psi(\mathbf{y}_l^*, \mathbf{y}_i^*). \end{aligned}$$

Then  $\widehat{T}_{k3} = n^{-2}(n-1)(n-2)[\widehat{T}_{k3}^* + (n-2)^{-1}\widehat{T}_{k1}^*]$ . By the triangle inequality, we deduce

$$(A.13) \quad \begin{aligned} |\widehat{T}_{k3} - T_{k3}| &= \left| \frac{(n-1)(n-2)}{n^2}(\widehat{T}_{k3}^* - T_{k3}) - \frac{3n-2}{n^2}T_{k3} \right. \\ &\quad \left. + \frac{n-1}{n^2}(\widehat{T}_{k1}^* - T_{k1}) + \frac{n-1}{n^2}T_{k1} \right| \leq |\widehat{T}_{k3}^* - T_{k3}| \\ &\quad + \left| \frac{3}{n}T_{k3} \right| + |\widehat{T}_{k1}^* - T_{k1}| + \left| \frac{1}{n}T_{k1} \right|. \end{aligned}$$

It follows from the Cauchy-Schwarz inequality and (17)–(18) that

$$\begin{aligned} T_{k3} &\leq \{E[\phi^2(X_{1k}^*, X_{2k}^*)]E[\psi^2(\mathbf{y}_1^*, \mathbf{y}_3^*)]\}^{1/2} \\ &\leq \{E[(X_{1k}^2 + X_{2k}^2)^2]E[(\|\tilde{\mathbf{y}}_1\|^2 + \|\tilde{\mathbf{y}}_3\|^2)^2]\}^{1/2} \\ &\leq \{E[2(X_{1k}^4 + X_{2k}^4)]E[2(\|\tilde{\mathbf{y}}_1\|^4 + \|\tilde{\mathbf{y}}_3\|^4)]\}^{1/2} \\ &= 4\{E(X_{1k}^4)E(\|\tilde{\mathbf{y}}_1\|^4)\}^{1/2} \end{aligned}$$

for all  $1 \leq k \leq p$ .

In Step 1, we have shown that  $T_{k1} \leq 4\{E(X_{1k}^4)E(\|\tilde{\mathbf{y}}_1\|^4)\}^{1/2}$  for all  $1 \leq k \leq p$ . By Condition 2,  $E(X_{1k}^4)$  and  $E(\|\mathbf{y}_1\|^4)$  are uniformly bounded from above by some positive constant for all  $1 \leq k \leq p$ . Note that  $T_{k1} \geq 0$  and  $T_{k3} \geq 0$ . Thus for any positive constant  $\tilde{C}$ , we have  $\max_{1 \leq k \leq p} |3n^{-1}T_{k3}| < \tilde{C}n^{-\kappa_2}/16$  and  $\max_{1 \leq k \leq p} |n^{-1}T_{k1}| < \tilde{C}n^{-\kappa_2}/16$  for all  $n$  sufficiently large. These two inequalities along with (A.13) entail

$$(A.14) \quad \begin{aligned} P(\max_{1 \leq k \leq p} |\widehat{T}_{k3} - T_{k3}| \geq \tilde{C}n^{-\kappa_2}/4) &\leq P(\max_{1 \leq k \leq p} |\widehat{T}_{k3}^* - T_{k3}| \\ &\geq \tilde{C}n^{-\kappa_2}/16) + P(\max_{1 \leq k \leq p} |\widehat{T}_{k1}^* - T_{k1}| \geq \tilde{C}n^{-\kappa_2}/16). \end{aligned}$$

Replacing  $\tilde{C}$  with  $\tilde{C}/2$  in (31) gives

$$(A.15) \quad \begin{aligned} P(\max_{1 \leq k \leq p} |\widehat{T}_{k1}^* - T_{k1}| \geq \tilde{C}n^{-\kappa_2}/16) &\leq p\tilde{C}_1 \exp\{-\tilde{C}_2 n^{(1-2\kappa_2)/3-2\eta}\} \\ &\quad + \tilde{C}_3 \exp\{-\tilde{C}_4 n^{3\eta/2}\}, \end{aligned}$$

where  $\tilde{C}_1, \dots, \tilde{C}_4$  are some positive constants.

It remains to bound  $P(\max_{1 \leq k \leq p} |\widehat{T}_{k3}^* - T_{k3}| \geq \tilde{C}n^{-\kappa_2}/16)$ . Let  $T_{k3} = T_{k3,1} + T_{k3,2} + T_{k3,3}$  with

$$\begin{aligned} T_{k3,1} &= E[\phi(X_{1k}^*, X_{2k}^*)\psi(\mathbf{y}_1^*, \mathbf{y}_3^*)\mathbb{I}\{\phi(X_{1k}^*, X_{2k}^*) \leq M_3\}\mathbb{I}\{\psi(\mathbf{y}_1^*, \mathbf{y}_3^*) \leq M_4\}], \\ T_{k3,2} &= E[\phi(X_{1k}^*, X_{2k}^*)\psi(\mathbf{y}_1^*, \mathbf{y}_3^*)\mathbb{I}\{\phi(X_{1k}^*, X_{2k}^*) \leq M_3\}\mathbb{I}\{\psi(\mathbf{y}_1^*, \mathbf{y}_3^*) > M_4\}], \\ T_{k3,3} &= E[\phi(X_{1k}^*, X_{2k}^*)\psi(\mathbf{y}_1^*, \mathbf{y}_3^*)\mathbb{I}\{\phi(X_{1k}^*, X_{2k}^*) > M_3\}]. \end{aligned}$$

Similarly, write  $\widehat{T}_{k3}^*$  as  $\widehat{T}_{k3}^* = \widehat{T}_{k3,1}^* + \widehat{T}_{k3,2}^* + \widehat{T}_{k3,3}^*$ , where

$$\begin{aligned}
\widehat{T}_{k3,1}^* &= \frac{1}{n(n-1)(n-2)} \sum_{i < j < l} [\phi(X_{ik}^*, X_{jk}^*) \psi(\mathbf{y}_i^*, \mathbf{y}_l^*) \mathbb{I}\{\phi(X_{ik}^*, X_{jk}^*) \leq M_3\} \\
&\quad \cdot \mathbb{I}\{\psi(\mathbf{y}_i^*, \mathbf{y}_l^*) \leq M_4\} \\
&\quad + \phi(X_{ik}^*, X_{lk}^*) \psi(\mathbf{y}_i^*, \mathbf{y}_j^*) \mathbb{I}\{\phi(X_{ik}^*, X_{lk}^*) \leq M_3\} \mathbb{I}\{\psi(\mathbf{y}_i^*, \mathbf{y}_j^*) \leq M_4\} \\
&\quad + \phi(X_{jk}^*, X_{ik}^*) \psi(\mathbf{y}_j^*, \mathbf{y}_l^*) \mathbb{I}\{\phi(X_{jk}^*, X_{ik}^*) \leq M_3\} \mathbb{I}\{\psi(\mathbf{y}_j^*, \mathbf{y}_l^*) \leq M_4\} \\
&\quad + \phi(X_{jk}^*, X_{lk}^*) \psi(\mathbf{y}_j^*, \mathbf{y}_i^*) \mathbb{I}\{\phi(X_{jk}^*, X_{lk}^*) \leq M_3\} \mathbb{I}\{\psi(\mathbf{y}_j^*, \mathbf{y}_i^*) \leq M_4\} \\
&\quad + \phi(X_{lk}^*, X_{ik}^*) \psi(\mathbf{y}_l^*, \mathbf{y}_j^*) \mathbb{I}\{\phi(X_{lk}^*, X_{ik}^*) \leq M_3\} \mathbb{I}\{\psi(\mathbf{y}_l^*, \mathbf{y}_j^*) \leq M_4\} \\
&\quad + \phi(X_{lk}^*, X_{jk}^*) \psi(\mathbf{y}_l^*, \mathbf{y}_i^*) \mathbb{I}\{\phi(X_{lk}^*, X_{jk}^*) \leq M_3\} \mathbb{I}\{\psi(\mathbf{y}_l^*, \mathbf{y}_i^*) \leq M_4\}] \\
&=: \frac{6}{n(n-1)(n-2)} \sum_{i < j < l} \widetilde{g}(X_{ik}^*, \mathbf{y}_i^*; X_{jk}^*, \mathbf{y}_j^*; X_{lk}^*, \mathbf{y}_l^*), \\
\widehat{T}_{k3,2}^* &= \frac{1}{n(n-1)(n-2)} \sum_{i < j < l} [\phi(X_{ik}^*, X_{jk}^*) \psi(\mathbf{y}_i^*, \mathbf{y}_l^*) \mathbb{I}\{\phi(X_{ik}^*, X_{jk}^*) \leq M_3\} \\
&\quad \cdot \mathbb{I}\{\psi(\mathbf{y}_i^*, \mathbf{y}_l^*) > M_4\} \\
&\quad + \phi(X_{ik}^*, X_{lk}^*) \psi(\mathbf{y}_i^*, \mathbf{y}_j^*) \mathbb{I}\{\phi(X_{ik}^*, X_{lk}^*) \leq M_3\} \mathbb{I}\{\psi(\mathbf{y}_i^*, \mathbf{y}_j^*) > M_4\} \\
&\quad + \phi(X_{jk}^*, X_{ik}^*) \psi(\mathbf{y}_j^*, \mathbf{y}_l^*) \mathbb{I}\{\phi(X_{jk}^*, X_{ik}^*) \leq M_3\} \mathbb{I}\{\psi(\mathbf{y}_j^*, \mathbf{y}_l^*) > M_4\} \\
&\quad + \phi(X_{jk}^*, X_{lk}^*) \psi(\mathbf{y}_j^*, \mathbf{y}_i^*) \mathbb{I}\{\phi(X_{jk}^*, X_{lk}^*) \leq M_3\} \mathbb{I}\{\psi(\mathbf{y}_j^*, \mathbf{y}_i^*) > M_4\} \\
&\quad + \phi(X_{lk}^*, X_{ik}^*) \psi(\mathbf{y}_l^*, \mathbf{y}_j^*) \mathbb{I}\{\phi(X_{lk}^*, X_{ik}^*) \leq M_3\} \mathbb{I}\{\psi(\mathbf{y}_l^*, \mathbf{y}_j^*) > M_4\} \\
&\quad + \phi(X_{lk}^*, X_{jk}^*) \psi(\mathbf{y}_l^*, \mathbf{y}_i^*) \mathbb{I}\{\phi(X_{lk}^*, X_{jk}^*) \leq M_3\} \mathbb{I}\{\psi(\mathbf{y}_l^*, \mathbf{y}_i^*) > M_4\}], \\
\widehat{T}_{k3,3}^* &= \frac{1}{n(n-1)(n-2)} \sum_{i < j < l} [\phi(X_{ik}^*, X_{jk}^*) \psi(\mathbf{y}_i^*, \mathbf{y}_l^*) \mathbb{I}\{\phi(X_{ik}^*, X_{jk}^*) > M_3\} \\
&\quad + \phi(X_{ik}^*, X_{lk}^*) \psi(\mathbf{y}_i^*, \mathbf{y}_j^*) \mathbb{I}\{\phi(X_{ik}^*, X_{lk}^*) > M_3\} \\
&\quad + \phi(X_{jk}^*, X_{ik}^*) \psi(\mathbf{y}_j^*, \mathbf{y}_l^*) \mathbb{I}\{\phi(X_{jk}^*, X_{ik}^*) > M_3\} \\
&\quad + \phi(X_{jk}^*, X_{lk}^*) \psi(\mathbf{y}_j^*, \mathbf{y}_i^*) \mathbb{I}\{\phi(X_{jk}^*, X_{lk}^*) > M_3\} \\
&\quad + \phi(X_{lk}^*, X_{ik}^*) \psi(\mathbf{y}_l^*, \mathbf{y}_j^*) \mathbb{I}\{\phi(X_{lk}^*, X_{ik}^*) > M_3\} \\
&\quad + \phi(X_{lk}^*, X_{jk}^*) \psi(\mathbf{y}_l^*, \mathbf{y}_i^*) \mathbb{I}\{\phi(X_{lk}^*, X_{jk}^*) > M_3\}].
\end{aligned}$$

Clearly,  $\widehat{T}_{k3,1}^*$ ,  $\widehat{T}_{k3,2}^*$ , and  $\widehat{T}_{k3,3}^*$  are unbiased estimators of  $T_{k3,1}$ ,  $T_{k3,2}$ , and

$T_{k3,3}$ , respectively. By the triangle inequality, we have

$$(A.16) \quad \begin{aligned} P(\max_{1 \leq k \leq p} |\widehat{T}_{k3}^* - T_{k3}| \geq \widetilde{C}n^{-\kappa_2}/16) \\ \leq \sum_{j=1}^3 P(\max_{1 \leq k \leq p} |\widehat{T}_{k3,j}^* - T_{k3,j}| \geq \widetilde{C}n^{-\kappa_2}/48). \end{aligned}$$

Note that  $\widetilde{g}$  defined in the expression for  $\widehat{T}_{k3,1}^*$  is the kernel of the  $U$ -statistic  $\widehat{T}_{k3,1}^*$  of order 3. Applying similar arguments to those for dealing with  $\widehat{T}_{k1,1}^*$  in Step 1 yields

$$\begin{aligned} P(|\widehat{T}_{k3,1}^* - T_{k3,1}| \geq \widetilde{C}n^{-\kappa_2}/48) &\leq 2 \exp\{-m_3 \widetilde{C}^2 n^{-2\kappa_2} / (1152 M_3^2 M_4^2)\} \\ &\leq 2 \exp\{-\widetilde{C}_5 n^{1-2\kappa_2-2\xi_3-2\xi_4}\} \end{aligned}$$

with  $\widetilde{C}_5$  some positive constant, by setting  $M_3 = n^{\xi_3}$  and  $M_4 = n^{\xi_4}$  with  $\xi_3, \xi_4 > 0$  and noting that  $m_3 = \lfloor n/3 \rfloor$ . Thus it follows from Bonferroni's inequality that

$$(A.17) \quad \begin{aligned} P(\max_{1 \leq k \leq p} |\widehat{T}_{k3,1}^* - T_{k3,1}| \geq \widetilde{C}n^{-\kappa_2}/48) \\ \leq \sum_{1 \leq k \leq p} P(|\widehat{T}_{k3,1}^* - T_{k3,1}| \geq \widetilde{C}n^{-\kappa_2}/48) \\ \leq 2p \exp\{-\widetilde{C}_5 n^{1-2\kappa_2-2\xi_3-2\xi_4}\}. \end{aligned}$$

Using similar arguments to those for (29)–(30), we can show that

$$(A.18) \quad \begin{aligned} P(\max_{1 \leq k \leq p} |\widehat{T}_{k3,2}^* - T_{k3,2}| \geq \widetilde{C}n^{-\kappa_2}/48) \\ \leq \widetilde{C}_6 n^{\kappa_2+\xi_3} \exp\{-2^{-3/2} c_0 n^{\xi_4/2}\}, \end{aligned}$$

$$(A.19) \quad \begin{aligned} P(\max_{1 \leq k \leq p} |\widehat{T}_{k3,3}^* - T_{k3,3}| \geq \widetilde{C}n^{-\kappa_2}/48) \\ \leq p \widetilde{C}_7 n^{\kappa_2} \exp(-8^{-1} c_0 n^{\xi_3}), \end{aligned}$$

where  $\widetilde{C}_6$  and  $\widetilde{C}_7$  are some positive constants.

Combining the results in (A.16)–(A.19) leads to

$$\begin{aligned} P(\max_{1 \leq k \leq p} |\widehat{T}_{k3}^* - T_{k3}| \geq \widetilde{C}n^{-\kappa_2}/16) &\leq 2p \exp\{-\widetilde{C}_5 n^{1-2\kappa_2-2\xi_3-2\xi_4}\} \\ &\quad + p \widetilde{C}_7 n^{\kappa_2} \exp(-8^{-1} c_0 n^{\xi_3}) + \widetilde{C}_6 n^{\kappa_2+\xi_3} \exp\{-2^{-3/2} c_0 n^{\xi_4/2}\}. \end{aligned}$$

Let  $\xi_3 = (1 - 2\kappa_2)/3 - 2\eta$  and  $\xi_4 = 3\eta$  with some  $0 < \eta < (1 - 2\kappa_2)/6$ . Then we have

$$P(\max_{1 \leq k \leq p} |\widehat{T}_{k3}^* - T_{k3}| \geq \widetilde{C}n^{-\kappa_2}/16) \leq p\widetilde{C}_8 \exp\{-\widetilde{C}_9n^{(1-2\kappa_2)/3-2\eta}\} \\ + \widetilde{C}_{10} \exp\{-\widetilde{C}_{11}n^{3\eta/2}\},$$

where  $\widetilde{C}_8, \dots, \widetilde{C}_{11}$  are some positive constants. This inequality together with (A.14)–(A.15) entails

$$(A.20) \quad P(\max_{1 \leq k \leq p} |\widehat{T}_{k3} - T_{k3}| \geq \widetilde{C}n^{-\kappa_2}/4) \leq p\widetilde{C}_1 \exp\{-\widetilde{C}_2n^{(1-2\kappa_2)/3-2\eta}\} \\ + \widetilde{C}_3 \exp\{-\widetilde{C}_4n^{3\eta/2}\}$$

for some positive constants  $\widetilde{C}_1, \dots, \widetilde{C}_4$ .

**D.4. Proof of Theorem 3.** For simplicity, we provide here only the proof for the case without variable screening, that is,  $d_1 = d_2 = p$ . The case with variable screening can be proved using similar arguments, in view of the sure screening property established in Theorem 1. By the definition of  $\widehat{\mathbf{B}}$ , we have

$$\frac{1}{2nq} \|\mathbf{Y} - \widetilde{\mathbf{X}}\widehat{\mathbf{B}}\|_F^2 + \lambda \|\widehat{\mathbf{B}}\|_{2,1} \leq \frac{1}{2nq} \|\mathbf{Y} - \widetilde{\mathbf{X}}\mathbf{B}^*\|_F^2 + \lambda \|\mathbf{B}^*\|_{2,1}.$$

Substituting  $Y = \widetilde{\mathbf{X}}\mathbf{B}^* + \mathbf{W}$  and rearranging terms yield

$$(A.21) \quad \frac{1}{2nq} \|\widetilde{\mathbf{X}}\widehat{\mathbf{\Delta}}\|_F^2 \leq \frac{1}{nq} \text{tr}(\mathbf{W}^T \widetilde{\mathbf{X}}\widehat{\mathbf{\Delta}}) + \lambda(\|\mathbf{B}^*\|_{2,1} - \|\widehat{\mathbf{B}}\|_{2,1}),$$

where  $\widehat{\mathbf{\Delta}} = \widehat{\mathbf{B}} - \mathbf{B}^*$ . An application of the Cauchy-Schwarz inequality gives

$$(A.22) \quad \text{tr}(\mathbf{W}^T \widetilde{\mathbf{X}}\widehat{\mathbf{\Delta}}) = \text{tr}(\widehat{\mathbf{\Delta}}\mathbf{W}^T \widetilde{\mathbf{X}}) = \sum_k \widehat{\mathbf{\Delta}}_k \left[ (\widetilde{\mathbf{X}}^T \mathbf{W})_k \right]^T \\ \leq \sum_k \|\widehat{\mathbf{\Delta}}_k\|_2 \cdot \|(\widetilde{\mathbf{X}}^T \mathbf{W})_k\|_2 \leq \|\widetilde{\mathbf{X}}^T \mathbf{W}\|_{2,\infty} \|\widehat{\mathbf{\Delta}}\|_{2,1},$$

where  $\widehat{\mathbf{\Delta}}_k$  and  $(\widetilde{\mathbf{X}}^T \mathbf{W})_k$  are the  $k$ th rows of  $\widehat{\mathbf{\Delta}}$  and  $\widetilde{\mathbf{X}}^T \mathbf{W}$ , respectively.

Note that  $q \leq p$  and  $\log p = o(n^\eta)$  with  $\eta = \min\{\eta_0, 1/2 - 2\xi\}$ . By Lemmas 7–8, with probability at least  $1 - O\{\exp(-\widetilde{C}_1n^{1/2-2\xi})\} - O(p^{-c}) = 1 - O(p^{-c_4})$  for some constants  $\widetilde{C}_1, c, c_4 > 0$  it holds that

$$(A.23) \quad \frac{1}{nq} \|\widetilde{\mathbf{X}}^T \mathbf{W}\|_{2,\infty} \leq \frac{\lambda}{2}$$

and

$$(A.24) \quad \min_{|J| \leq s, \mathbf{\Delta} \in \mathbb{R}^{\bar{p} \times q} \setminus \{\mathbf{0}\}, \|\mathbf{\Delta}_{J^c}\|_{2,1} \leq 3\|\mathbf{\Delta}_J\|_{2,1}} \frac{\|\tilde{\mathbf{X}}\mathbf{\Delta}\|_F}{\sqrt{n}\|\mathbf{\Delta}_J\|_F} \geq \frac{\kappa}{2}.$$

From now on, we condition on the event that these two inequalities hold. In view of (A.21) and (A.22), we have the following basic inequality

$$(A.25) \quad \frac{1}{2nq} \|\tilde{\mathbf{X}}\hat{\mathbf{\Delta}}\|_F^2 + \frac{\lambda}{2} \|\hat{\mathbf{B}} - \mathbf{B}^*\|_{2,1} \leq \lambda(\|\mathbf{B}^*\|_{2,1} - \|\hat{\mathbf{B}}\|_{2,1} + \|\hat{\mathbf{B}} - \mathbf{B}^*\|_{2,1}) \leq 2\lambda\|(\hat{\mathbf{B}} - \mathbf{B}^*)_{S^c}\|_{2,1},$$

where we have used the fact that  $\|(\mathbf{B}^*)_{S^c}\|_{2,1} - \|(\hat{\mathbf{B}})_{S^c}\|_{2,1} + \|(\hat{\mathbf{B}} - \mathbf{B}^*)_{S^c}\|_{2,1} = 0$  and the triangle inequality.

The basic inequality (A.25) implies

$$(A.26) \quad \frac{1}{2nq} \|\tilde{\mathbf{X}}\hat{\mathbf{\Delta}}\|_F^2 \leq 2\lambda\|\hat{\mathbf{\Delta}}_S\|_{2,1} \leq 2\lambda\sqrt{s}\|\hat{\mathbf{\Delta}}_S\|_F,$$

where the last inequality holds since

$$(A.27) \quad \|\hat{\mathbf{\Delta}}_S\|_{2,1} = \sum_{k \in S} \|\hat{\mathbf{\Delta}}_k\|_2 \leq \sqrt{s \sum_{k \in S} \|\hat{\mathbf{\Delta}}_k\|_2^2} = \sqrt{s}\|\hat{\mathbf{\Delta}}_S\|_F.$$

Moreover, it follows from (A.25) that  $\|\hat{\mathbf{\Delta}}_{S^c}\|_{2,1} \leq 3\|\hat{\mathbf{\Delta}}_S\|_{2,1}$  and thus by (A.24), we have  $\|\hat{\mathbf{\Delta}}_S\|_F \leq 2\|\tilde{\mathbf{X}}\hat{\mathbf{\Delta}}\|_F/(\kappa\sqrt{n})$ . This inequality along with (A.25)–(A.27) yields

$$\frac{1}{2nq} \|\tilde{\mathbf{X}}\hat{\mathbf{\Delta}}\|_F^2 + \frac{\lambda}{2} \|\hat{\mathbf{B}} - \mathbf{B}^*\|_{2,1} \leq \frac{4\lambda\sqrt{s}\|\tilde{\mathbf{X}}\hat{\mathbf{\Delta}}\|_F}{\kappa\sqrt{n}},$$

which gives  $n^{-1/2}\|\tilde{\mathbf{X}}\hat{\mathbf{\Delta}}\|_F \leq 8q\lambda\sqrt{s}/\kappa$ . Therefore, we obtain

$$\frac{1}{2nq} \|\tilde{\mathbf{X}}(\hat{\mathbf{B}} - \mathbf{B}^*)\|_F^2 + \frac{\lambda}{2} \|\hat{\mathbf{B}} - \mathbf{B}^*\|_{2,1} \leq \frac{32qs\lambda^2}{\kappa^2},$$

which completes the proof for the first part of Theorem 3.

We next proceed to prove the second part of Theorem 3. We condition on the event that (A.4) holds. Denote by  $S(\mathbf{B})$  the row support of any matrix  $\mathbf{B}$ . We need to show that with the same probability  $S(\hat{\mathbf{B}}) = S(\mathbf{B}^*)$  holds. To this end, we first prove  $S(\mathbf{B}^*) \subset S(\hat{\mathbf{B}})$ . For any  $j_0 \in S(\mathbf{B}^*)$ , if  $j_0 \notin S(\hat{\mathbf{B}})$  then the  $j_0$ th row of  $\hat{\mathbf{B}}$  is zero, which means  $\|\hat{\mathbf{B}}_{j_0}\| \leq 64c_3\kappa^{-2}s\sqrt{q(\log p)/n}$ .



It then follows from the condition of  $\min_{j \in S} \|\mathbf{B}_j^*\| > 128c_3\kappa^{-2}s\sqrt{q(\log p)/n}$  that

$$\begin{aligned} \frac{1}{\sqrt{q}} \|\widehat{\mathbf{B}} - \mathbf{B}^*\|_{2,1} &\geq \frac{1}{\sqrt{q}} \|\widehat{\mathbf{B}}_{j_0} - \mathbf{B}_{j_0}^*\| \geq \frac{1}{\sqrt{q}} (\|\mathbf{B}_{j_0}^*\| - \|\widehat{\mathbf{B}}_{j_0}\|) \\ &> 64c_3\kappa^{-2}s\sqrt{(\log p)/n}, \end{aligned}$$

which leads to a contradiction to the estimation bound (A.4). Thus it holds that  $S(\mathbf{B}^*) \subset S(\widetilde{\mathbf{B}})$ . We can also show that  $S(\widetilde{\mathbf{B}}) \subset S(\mathbf{B}^*)$ . In fact, if there exists some  $j_0$  such that  $j_0 \in S(\widetilde{\mathbf{B}})$  and  $j_0 \notin S(\mathbf{B}^*)$ , then we have  $\|\widehat{\mathbf{B}}_{j_0}\| > 64c_3\kappa^{-2}s\sqrt{q(\log p)/n}$  and  $\mathbf{B}_{j_0}^* = \mathbf{0}$ , and thus

$$\frac{1}{\sqrt{q}} \|\widehat{\mathbf{B}} - \mathbf{B}^*\|_{2,1} \geq \frac{1}{\sqrt{q}} \|\widehat{\mathbf{B}}_{j_0} - \mathbf{B}_{j_0}^*\| > 64c_3\kappa^{-2}s\sqrt{(\log p)/n},$$

which contradicts again the bound (A.4). Combining these results yields that with the same probability, the row support of  $\widetilde{\mathbf{B}}$  is identical to the true row support  $S$ .

We finally prove (A.5). By assumption, the RE( $2s$ ) condition holds. Using similar arguments as for proving (A.23)–(A.24), we can show that with probability at least  $1 - O\{\exp(-\widetilde{C}_1 \cdot n^{1/2-2\xi})\} - O(p^{-c}) = 1 - O(p^{-c_4})$  for some constants  $\widetilde{C}_1, c, c_4 > 0$ , it holds that

$$(A.28) \quad \frac{1}{nq} \|\widetilde{\mathbf{X}}^T \mathbf{W}\|_{2,\infty} \leq \frac{\lambda}{2}$$

and

$$(A.29) \quad \min_{|J| \leq 2s, \mathbf{\Delta} \in \mathbb{R}^{\widetilde{p} \times q} \setminus \{\mathbf{0}\}, \|\mathbf{\Delta}_{J^c}\|_{2,1} \leq 3\|\mathbf{\Delta}_J\|_{2,1}} \frac{\|\widetilde{\mathbf{X}}\mathbf{\Delta}\|_F}{\sqrt{n}\|\mathbf{\Delta}_J\|_F} \geq \frac{\kappa(2s)}{2}.$$

Recall that  $\widehat{\mathbf{\Delta}} = \widehat{\mathbf{B}} - \mathbf{B}^*$ . Let  $S'$  be a subset of  $S^c$  corresponding to the  $s$  largest values of  $\|\widehat{\mathbf{\Delta}}_k\|$ . Then we have  $|S \cup S'| = 2s$ . From now on, we condition on the event that inequalities (A.28) and (A.29) hold. Conditional on such an event, the basic inequality (A.25) still holds. Thus we have  $\|\widehat{\mathbf{\Delta}}_{S^c}\|_{2,1} \leq 3\|\widehat{\mathbf{\Delta}}_S\|_{2,1}$ , which entails

$$\|\widehat{\mathbf{\Delta}}_{(S \cup S')^c}\|_{2,1} \leq \|\widehat{\mathbf{\Delta}}_{S^c}\|_{2,1} \leq 3\|\widehat{\mathbf{\Delta}}_S\|_{2,1} \leq 3\|\widehat{\mathbf{\Delta}}_{S \cup S'}\|_{2,1}.$$

This together with (A.29) yields  $\|\widehat{\mathbf{\Delta}}_{S \cup S'}\|_F \leq 2\|\widetilde{\mathbf{X}}\widehat{\mathbf{\Delta}}\|_F / (\kappa(2s)\sqrt{n})$ . From (A.26), we have

$$\frac{1}{2nq} \|\widetilde{\mathbf{X}}\widehat{\mathbf{\Delta}}\|_F^2 \leq 2\lambda\sqrt{s}\|\widehat{\mathbf{\Delta}}_S\|_F \leq 2\lambda\sqrt{s}\|\widehat{\mathbf{\Delta}}_{S \cup S'}\|_F.$$

Combining these two results gives

$$(A.30) \quad \|\widehat{\Delta}_{S \cup S'}\|_F \leq 16q\lambda\sqrt{s}/\kappa^2(2s).$$

Since the  $j$ th largest norm in the set  $\{\|\widehat{\Delta}_k\| : k \in S^c\}$  is bounded from above by  $\|\widehat{\Delta}_{S^c}\|_{2,1}/j$ , it holds that

$$\begin{aligned} \sum_{k \in (S \cup S')^c} \|\widehat{\Delta}_k\|^2 &\leq \sum_{k=s+1}^{\bar{p}-s} \frac{\|\widehat{\Delta}_{S^c}\|_{2,1}^2}{k^2} \leq \frac{\|\widehat{\Delta}_{S^c}\|_{2,1}^2}{s} \leq \frac{9\|\widehat{\Delta}_S\|_{2,1}^2}{s} \\ &\leq 9 \sum_{k \in S} \|\widehat{\Delta}_k\|^2 \leq 9 \sum_{k \in S \cup S'} \|\widehat{\Delta}_k\|^2, \end{aligned}$$

which results in  $\|\widehat{\Delta}\|_F^2 \leq 10\|\widehat{\Delta}_{S \cup S'}\|_F^2$ . This inequality along with (A.30) yields

$$\frac{1}{\sqrt{q}}\|\widehat{\Delta}\|_F \leq \frac{16\sqrt{10}c_3}{\kappa^2(2s)}\sqrt{s(\log p)/n},$$

which concludes the proof for the third part of Theorem 3.

## APPENDIX E: ADDITIONAL TECHNICAL DETAILS AND LEMMAS

**E.1. Terms  $E(Y|X_j)$  and  $E(Y^2|X_j)$  under model (2).** Since the covariates  $X_1, \dots, X_p$  are all independent with mean zero and the random error  $W$  is of mean zero and independent of all  $X_j$ 's, it is immediate that  $E(Y|X_j) = \alpha + \beta_j X_j$ . We now calculate  $E(Y^2|X_j)$ . Define

$$\begin{aligned} J_1 &= \sum_{j=1}^p \beta_j X_j, & J_2 &= \sum_{k=1}^{p-1} \sum_{\ell=k+1}^p \gamma_{k\ell} X_k X_\ell, & J_3 &= \sum_{k \neq j} \beta_k X_k, \\ J_4 &= \sum_{k=1}^{j-1} \gamma_{kj} X_k + \sum_{\ell=j+1}^p \gamma_{j\ell} X_\ell, & J_5 &= \sum_{k=1, k \neq j}^{p-1} \sum_{\ell=k+1, \ell \neq j}^p \gamma_{k\ell} X_k X_\ell. \end{aligned}$$

Then we have  $Y = \alpha + J_1 + J_2 + W$  with  $J_1 = \beta_j X_j + J_3$  and  $J_2 = J_4 X_j + J_5$ , and  $J_3, J_4$ , and  $J_5$  are independent of  $X_j$ . Applying the properties of conditional expectation yields

$$\begin{aligned} E[(\alpha + J_1 + J_2)W|X_j] &= E\{E[(\alpha + J_1 + J_2)W|X_1, \dots, X_p]|X_j\} \\ &= E[(\alpha + J_1 + J_2)E(W|X_1, \dots, X_p)|X_j] = 0 \end{aligned}$$

and

$$\begin{aligned}
E[(\alpha + J_1 + J_2)^2 | X_j] &= E\{[(\beta_j + J_4)X_j + (\alpha + J_3 + J_5)]^2 | X_j\} \\
&= X_j^2 E[(\beta_j + J_4)^2] + 2X_j E[(\beta_j + J_4)(\alpha + J_3 + J_5)] + E[(\alpha + J_3 + J_5)^2] \\
&= \left[ \beta_j^2 + \sum_{k=1}^{j-1} \gamma_{kj}^2 E(X_k^2) + \sum_{\ell=j+1}^p \gamma_{j\ell}^2 E(X_\ell^2) \right] X_j^2 \\
&\quad + 2 \left[ \beta_j \alpha + \sum_{k=1}^{j-1} \beta_k \gamma_{kj} E(X_k^2) + \sum_{\ell=j+1}^p \beta_\ell \gamma_{j\ell} E(X_\ell^2) \right] X_j \\
&\quad + \alpha^2 + \sum_{k \neq j} \beta_k^2 E(X_k^2) + \sum_{k=1, k \neq j}^{p-1} \sum_{\ell=k+1, \ell \neq j}^p \gamma_{k\ell}^2 E(X_k^2) E(X_\ell^2).
\end{aligned}$$

Therefore, it holds that

$$\begin{aligned}
E(Y^2 | X_j) &= E[(\alpha + J_1 + J_2)^2 | X_j] + 2E[(\alpha + J_1 + J_2)W | X_j] + E(W^2 | X_j) \\
&= \left[ \beta_j^2 + \sum_{k=1}^{j-1} \gamma_{kj}^2 E(X_k^2) + \sum_{\ell=j+1}^p \gamma_{j\ell}^2 E(X_\ell^2) \right] X_j^2 \\
&\quad + 2 \left[ \beta_j \alpha + \sum_{k=1}^{j-1} \beta_k \gamma_{kj} E(X_k^2) + \sum_{\ell=j+1}^p \beta_\ell \gamma_{j\ell} E(X_\ell^2) \right] X_j + C_j,
\end{aligned}$$

where  $C_j = \alpha^2 + \sum_{k \neq j} \beta_k^2 E(X_k^2) + \sum_{k=1, k \neq j}^{p-1} \sum_{\ell=k+1, \ell \neq j}^p \gamma_{k\ell}^2 E(X_k^2) E(X_\ell^2) + \sigma^2$  is a constant that is free of  $X_j$ , and  $\sigma^2$  is the variance of  $W$ .

## E.2. Lemma 1 and its proof.

LEMMA 1. *Let  $\hat{A}$  and  $\hat{B}$  be estimates of  $A$  and  $B$ , respectively, based on a sample of size  $n$ . Assume that both  $A$  and  $B$  are bounded and for any constant  $\tilde{C} > 0$ , there exist positive constants  $\tilde{C}_1, \dots, \tilde{C}_4$  such that*

$$\begin{aligned}
P\left(|\hat{A} - A| \geq \tilde{C}n^{-\kappa}\right) &\leq \tilde{C}_1 \exp\left\{-\tilde{C}_2 n^{f(\kappa)}\right\} \\
P\left(|\hat{B} - B| \geq \tilde{C}n^{-\kappa}\right) &\leq \tilde{C}_3 \exp\left\{-\tilde{C}_4 n^{f(\kappa)}\right\}
\end{aligned}$$

with  $f(\kappa)$  some function of  $\kappa$ . Then for any constant  $\tilde{C} > 0$ , there exist positive constants  $\tilde{C}_5$  and  $\tilde{C}_6$  such that

$$P(|\hat{A}\hat{B} - AB| \geq \tilde{C}n^{-\kappa}) \leq \tilde{C}_5 \exp\left\{-\tilde{C}_6 n^{f(\kappa)}\right\}.$$

*Proof.* Note that  $|\widehat{A}\widehat{B} - AB| \leq |\widehat{A}(\widehat{B} - B)| + |(\widehat{A} - A)B|$ . Thus for any positive constant  $\widetilde{C}$ , we have

$$(A.31) \quad P(|\widehat{A}\widehat{B} - AB| \geq \widetilde{C}n^{-\kappa}) \leq P(|\widehat{A}(\widehat{B} - B)| \geq \widetilde{C}n^{-\kappa}/2) \\ + P(|(\widehat{A} - A)B| \geq \widetilde{C}n^{-\kappa}/2).$$

We first deal with the second term on the right hand side of (A.31). Since both  $A$  and  $B$  are bounded, there exists some positive constant  $L$  such that  $|A| \leq L$  and  $|B| \leq L$ . It follows that

$$(A.32) \quad P(|(\widehat{A} - A)B| \geq \widetilde{C}n^{-\kappa}/2) \leq P(|\widehat{A} - A|L \geq \widetilde{C}n^{-\kappa}/2) \\ = P\{|\widehat{A} - A| \geq (2L)^{-1}\widetilde{C}n^{-\kappa}\} \leq \widetilde{C}_1 \exp\left\{-\widetilde{C}_2 n^{f(\kappa_1)}\right\},$$

where  $\widetilde{C}_1$  and  $\widetilde{C}_2$  are some positive constants.

We next consider the first term on the right hand side of (A.31). Note that

$$(A.33) \quad P(|\widehat{A}(\widehat{B} - B)| \geq \widetilde{C}n^{-\kappa}/2) \leq P\left\{|\widehat{A}(\widehat{B} - B)| \geq \widetilde{C}n^{-\kappa}/2, \right. \\ \left. |\widehat{A}| \geq L + \frac{\widetilde{C}}{2}n^{-\kappa}\right\} + P\left(|\widehat{A}(\widehat{B} - B)| \geq \frac{\widetilde{C}}{2}n^{-\kappa}, |\widehat{A}| < L + \frac{\widetilde{C}}{2}n^{-\kappa}\right) \\ \leq P(|\widehat{A}| \geq L + \frac{\widetilde{C}}{2}n^{-\kappa}) + P(|\widehat{A}(\widehat{B} - B)| \geq \frac{\widetilde{C}}{2}n^{-\kappa}, |\widehat{A}| < L + \widetilde{C}) \\ \leq P(|\widehat{A}| \geq L + \widetilde{C}n^{-\kappa}/2) + P\{(L + \widetilde{C})|\widehat{B} - B| \geq \widetilde{C}n^{-\kappa}/2\}.$$

We will bound the two terms on the right hand side of (A.33) separately. It follows from  $|A| \leq L$  that

$$(A.34) \quad P(|\widehat{A}| \geq L + \widetilde{C}n^{-\kappa}/2) \leq P(|\widehat{A} - A| + |A| \geq L + \widetilde{C}n^{-\kappa}/2) \\ \leq P\{|\widehat{A} - A| \geq 2^{-1}\widetilde{C}n^{-\kappa}\} \leq \widetilde{C}_3 \exp\left\{-\widetilde{C}_4 n^{f(\kappa_1)}\right\},$$

where  $\widetilde{C}_3$  and  $\widetilde{C}_4$  are some positive constants. It also holds that

$$P((L + \widetilde{C})|\widehat{B} - B| \geq \widetilde{C}n^{-\kappa}/2) = P\{|\widehat{B} - B| \geq (2L + 2\widetilde{C})^{-1}\widetilde{C}n^{-\kappa}\} \\ \leq \widetilde{C}_7 \exp\left\{-\widetilde{C}_8 n^{f(\kappa)}\right\},$$

where  $\widetilde{C}_7$  and  $\widetilde{C}_8$  are some positive constants. This inequality together with (A.31)–(A.34) entails

$$P(|\widehat{A}\widehat{B} - AB| \geq \widetilde{C}n^{-\kappa}) \leq \widetilde{C}_1 \exp\left\{-\widetilde{C}_2 n^{f(\kappa)}\right\} + \widetilde{C}_3 \exp\left\{-\widetilde{C}_4 n^{f(\kappa)}\right\} \\ + \widetilde{C}_7 \exp\left\{-\widetilde{C}_8 n^{f(\kappa)}\right\} \leq \widetilde{C}_5 \exp\left\{-\widetilde{C}_6 n^{f(\kappa)}\right\},$$

where  $\widetilde{C}_5 = \widetilde{C}_1 + \widetilde{C}_3 + \widetilde{C}_7$  and  $\widetilde{C}_6 = \min\{\widetilde{C}_2, \widetilde{C}_4, \widetilde{C}_8\}$ .

**E.3. Lemma 2 and its proof.** For any set  $\mathcal{D}$ , we denote by  $|\mathcal{D}|$  its cardinality throughout the paper.

LEMMA 2. Let  $\widehat{B}_j \geq 0$  be an estimate of  $B_j$  based on a sample of size  $n$  for each  $j \in \mathcal{D} \subset \{1, \dots, p\}$ . Assume that  $\min_{j \in \mathcal{D}} B_j \geq L$  for some positive constant  $L$ , and for any constant  $\widetilde{C} > 0$ , there exist positive constants  $\widetilde{C}_1$  and  $\widetilde{C}_2$  such that

$$P\left(\max_{j \in \mathcal{D}} |\widehat{B}_j - B_j| \geq \widetilde{C}n^{-\kappa}\right) \leq |\mathcal{D}| \widetilde{C}_1 \exp\left\{-\widetilde{C}_2 n^{f(\kappa)}\right\}$$

with  $f(\kappa)$  some function of  $\kappa$ . Then for any constant  $\widetilde{C} > 0$ , there exist positive constants  $\widetilde{C}_3$  and  $\widetilde{C}_4$  such that

$$P\left(\max_{j \in \mathcal{D}} |\sqrt{\widehat{B}_j} - \sqrt{B_j}| \geq \widetilde{C}n^{-\kappa}\right) \leq |\mathcal{D}| \widetilde{C}_3 \exp\left\{-\widetilde{C}_4 n^{f(\kappa)}\right\}.$$

*Proof.* Since  $\min_{j \in \mathcal{D}} B_j \geq L$  for some positive constant  $L$ , there exists a constant  $L_0$  such that  $0 < L_0 < L$ . Note that for any positive constant  $\widetilde{C}$ ,

$$\begin{aligned} \text{(A.35)} \quad P(\max_{j \in \mathcal{D}} |\sqrt{\widehat{B}_j} - \sqrt{B_j}| \geq \widetilde{C}n^{-\kappa}) &\leq P\left\{\max_{j \in \mathcal{D}} |\sqrt{\widehat{B}_j} - \sqrt{B_j}| \geq \widetilde{C}n^{-\kappa}, \right. \\ &\quad \left. \min_{j \in \mathcal{D}} |\widehat{B}_j| \leq L - L_0 n^{-\kappa}\right\} + P\left\{\max_{j \in \mathcal{D}} |\sqrt{\widehat{B}_j} - \sqrt{B_j}| \geq \widetilde{C}n^{-\kappa}, \right. \\ &\quad \left. \min_{j \in \mathcal{D}} |\widehat{B}_j| > L - L_0 n^{-\kappa}\right\} \leq P(\min_{j \in \mathcal{D}} |\widehat{B}_j| \leq L - L_0 n^{-\kappa}) \\ &\quad + P(\max_{j \in \mathcal{D}} \frac{|\widehat{B}_j - B_j|}{|\sqrt{\widehat{B}_j} + \sqrt{B_j}|} \geq \widetilde{C}n^{-\kappa}, \min_{j \in \mathcal{D}} |\widehat{B}_j| > L - L_0). \end{aligned}$$

We first consider the first term on the right hand side of (A.35). It follows from  $\min_{j \in \mathcal{D}} B_j \geq L$  that

$$\begin{aligned} \text{(A.36)} \quad P(\min_{j \in \mathcal{D}} |\widehat{B}_j| \leq L - L_0 n^{-\kappa}) &\leq P\left\{\min_{j \in \mathcal{D}} |B_j| - \max_{j \in \mathcal{D}} |\widehat{B}_j - B_j| \right. \\ &\quad \left. \leq L - L_0 n^{-\kappa}\right\} \leq P(\max_{j \in \mathcal{D}} |\widehat{B}_j - B_j| \geq L_0 n^{-\kappa}) \\ &\leq |\mathcal{D}| \widetilde{C}_1 \exp\left\{-\widetilde{C}_2 n^{f(\kappa)}\right\}, \end{aligned}$$

where  $\widetilde{C}_1$  and  $\widetilde{C}_2$  are some positive constants.

Next we consider the second term on the right hand side of (A.35). For any positive constant  $\tilde{C}$ , we have

$$\begin{aligned}
\text{(A.37)} \quad & P\left(\max_{j \in \mathcal{D}} \frac{|\hat{B}_j - B_j|}{|\sqrt{\hat{B}_j} + \sqrt{B_j}|} \geq \tilde{C}n^{-\kappa}, \min_{j \in \mathcal{D}} |\hat{B}_j| > L - L_0\right) \\
& \leq P\left\{\max_{j \in \mathcal{D}} |\hat{B}_j - B_j| \geq \tilde{C}(\sqrt{L - L_0} + \sqrt{L})n^{-\kappa}\right\} \\
& \leq |\mathcal{D}|\tilde{C}_5 \exp\left\{-\tilde{C}_6 n^{f(\kappa)}\right\},
\end{aligned}$$

where  $\tilde{C}_5$  and  $\tilde{C}_6$  are some positive constants. Combining (A.35)–(A.37) yields

$$\text{(A.38)} \quad P\left(\max_{j \in \mathcal{D}} |\sqrt{\hat{B}_j} - \sqrt{B_j}| \geq \tilde{C}n^{-\kappa}\right) \leq |\mathcal{D}|\tilde{C}_3 \exp\left\{-\tilde{C}_4 n^{f(\kappa)}\right\},$$

where  $\tilde{C}_3 = \tilde{C}_1 + \tilde{C}_5$  and  $\tilde{C}_4 = \min\{\tilde{C}_2, \tilde{C}_6\}$ .

#### E.4. Lemma 3 and its proof.

LEMMA 3. *Let  $\hat{A}_j$  and  $\hat{B}_j$  be estimates of  $A_j$  and  $B_j$ , respectively, based on a sample of size  $n$  for each  $j \in \mathcal{D} \subset \{1, \dots, p\}$ . Assume that  $A_j$  and  $B_j$  satisfy  $\max_{j \in \mathcal{D}} |A_j| \leq L_1$  and  $\min_{j \in \mathcal{D}} |B_j| \geq L_2$  for some constants  $L_1, L_2 > 0$ , and for any constant  $\tilde{C} > 0$ , there exist constants  $\tilde{C}_1, \dots, \tilde{C}_6 > 0$  such that*

$$\begin{aligned}
P\left(\max_{j \in \mathcal{D}} |\hat{A}_j - A_j| \geq \tilde{C}n^{-\kappa}\right) & \leq |\mathcal{D}|\tilde{C}_1 \exp\left\{-\tilde{C}_2 n^{f(\kappa)}\right\} + \tilde{C}_3 \exp\left\{-\tilde{C}_4 n^{g(\kappa)}\right\}, \\
P\left(\max_{j \in \mathcal{D}} |\hat{B}_j - B_j| \geq \tilde{C}n^{-\kappa}\right) & \leq |\mathcal{D}|\tilde{C}_5 \exp\left\{-\tilde{C}_6 n^{f(\kappa)}\right\}
\end{aligned}$$

with  $f(\kappa)$  and  $g(\kappa)$  some functions of  $\kappa$ . Then for any constant  $\tilde{C} > 0$ , there exist positive constants  $\tilde{C}_7, \dots, \tilde{C}_{10}$  such that

$$\begin{aligned}
P\left(\max_{j \in \mathcal{D}} \left|\frac{\hat{A}_j}{\hat{B}_j} - \frac{A_j}{B_j}\right| \geq \tilde{C}n^{-\kappa}\right) & \leq |\mathcal{D}|\tilde{C}_7 \exp\left\{-\tilde{C}_8 n^{f(\kappa)}\right\} \\
& \quad + \tilde{C}_9 \exp\left\{-\tilde{C}_{10} n^{g(\kappa)}\right\}.
\end{aligned}$$

*Proof.* Since  $\min_{j \in \mathcal{D}} |B_j| \geq L_2 > 0$ , there exists some constant  $L_0$  such

that  $0 < L_0 < L_2$ . Note that for any positive constant  $\tilde{C}$ , we have

$$\begin{aligned}
\text{(A.39)} \quad P(\max_{j \in \mathcal{D}} |\frac{\hat{A}_j}{\hat{B}_j} - \frac{A_j}{B_j}| \geq \tilde{C}n^{-\kappa}) &\leq P\left\{ \max_{j \in \mathcal{D}} |\frac{\hat{A}_j}{\hat{B}_j} - \frac{A_j}{B_j}| \geq \tilde{C}n^{-\kappa}, \right. \\
&\quad \left. \min_{j \in \mathcal{D}} |\hat{B}_j| \leq L_2 - L_0n^{-\kappa} \right\} + P\left\{ \max_{j \in \mathcal{D}} |\frac{\hat{A}_j}{\hat{B}_j} - \frac{A_j}{B_j}| \geq \tilde{C}n^{-\kappa}, \right. \\
&\quad \left. \min_{j \in \mathcal{D}} |\hat{B}_j| > L_2 - L_0n^{-\kappa} \right\} \leq P(\min_{j \in \mathcal{D}} |\hat{B}_j| \leq L_2 - L_0n^{-\kappa}) \\
&\quad + P(\max_{j \in \mathcal{D}} |\frac{\hat{A}_j}{\hat{B}_j} - \frac{A_j}{B_j}| \geq \tilde{C}n^{-\kappa}, \min_{j \in \mathcal{D}} |\hat{B}_j| > L_2 - L_0).
\end{aligned}$$

We start with the first term on the right hand side of (A.39). In light of  $\min_{j \in \mathcal{D}} |B_j| \geq L_2$ , we deduce

$$\begin{aligned}
\text{(A.40)} \quad P(\min_{j \in \mathcal{D}} |\hat{B}_j| \leq L_2 - L_0n^{-\kappa}) &\leq P\left\{ \min_{j \in \mathcal{D}} |B_j| - \max_{j \in \mathcal{D}} |\hat{B}_j - B_j| \right. \\
&\quad \left. \leq L_2 - L_0n^{-\kappa} \right\} \leq P(\max_{j \in \mathcal{D}} |\hat{B}_j - B_j| \geq L_0n^{-\kappa}) \\
&\leq |\mathcal{D}| \tilde{C}_1 \exp\left\{-\tilde{C}_2 n^{f(\kappa)}\right\},
\end{aligned}$$

where  $\tilde{C}_1$  and  $\tilde{C}_2$  are some positive constants.

The second term on the right hand side of (A.39) can be bounded as

$$\begin{aligned}
\text{(A.41)} \quad P(\max_{j \in \mathcal{D}} |\frac{\hat{A}_j}{\hat{B}_j} - \frac{A_j}{B_j}| \geq \tilde{C}n^{-\kappa}, \min_{j \in \mathcal{D}} |\hat{B}_j| > L_2 - L_0) \\
\leq P(\max_{j \in \mathcal{D}} |\frac{\hat{A}_j}{\hat{B}_j} - \frac{A_j}{B_j}| \geq \tilde{C}n^{-\kappa}/2, \min_{j \in \mathcal{D}} |\hat{B}_j| > L_2 - L_0) \\
\quad + P(\max_{j \in \mathcal{D}} |\frac{A_j}{\hat{B}_j} - \frac{A_j}{B_j}| \geq \tilde{C}n^{-\kappa}/2, \min_{j \in \mathcal{D}} |\hat{B}_j| > L_2 - L_0) \\
\leq P\{\max_{j \in \mathcal{D}} |\hat{A}_j - A_j| \geq 2^{-1}(L_2 - L_0)\tilde{C}n^{-\kappa}\} \\
\quad + P\{\max_{j \in \mathcal{D}} |\hat{B}_j - B_j| \geq (2L_1)^{-1}(L_2 - L_0)L_2\tilde{C}n^{-\kappa}\} \\
\leq |\mathcal{D}| \tilde{C}_3 \exp\left\{-\tilde{C}_4 n^{f(\kappa)}\right\} + \tilde{C}_9 \exp\left\{-\tilde{C}_{10} n^{g(\kappa)}\right\} \\
\quad + |\mathcal{D}| \tilde{C}_5 \exp\left\{-\tilde{C}_6 n^{f(\kappa)}\right\},
\end{aligned}$$

where  $\tilde{C}_3, \dots, \tilde{C}_6$ , and  $\tilde{C}_9, \tilde{C}_{10}$  are some positive constants. Combining (A.39)–

(A.41) results in

$$P(\max_{j \in \mathcal{D}} |\frac{\hat{A}_j}{\hat{B}_j} - \frac{A_j}{B_j}| \geq \tilde{C}n^{-\kappa}) \leq |\mathcal{D}|\tilde{C}_7 \exp\{-\tilde{C}_8 n^{f(\kappa)}\} \\ + \tilde{C}_9 \exp\{-\tilde{C}_{10} n^{g(\kappa)}\},$$

where  $\tilde{C}_7 = \tilde{C}_1 + \tilde{C}_3 + \tilde{C}_5$  and  $\tilde{C}_8 = \min\{\tilde{C}_2, \tilde{C}_4, \tilde{C}_6\}$ .

### E.5. Lemma 4 and its proof.

LEMMA 4. *Let  $Z$  be a nonnegative random variable satisfying  $P(Z > t) \leq \tilde{C}_1 \exp(-\tilde{C}_2 t^2)$  for all  $t > 0$  with  $\tilde{C}_1, \tilde{C}_2 > 0$  some constants. Then*

$$E \left[ \exp \left( \frac{\tilde{C}_2}{2} Z^2 \right) \right] \leq 1 + \tilde{C}_1 \quad \text{and} \quad E(Z^{2m}) \leq (1 + \tilde{C}_1)(2\tilde{C}_2^{-1})^m m!$$

for any nonnegative integer  $m$ .

*Proof.* Let  $F(t)$  be the cumulative distribution function of  $Z$ . Then

$$1 - F(t) = P(Z > t) \leq \tilde{C}_1 \exp(-\tilde{C}_2 t^2)$$

for all  $t > 0$ . Using integration by parts, we have

$$E \left[ \exp \left( \frac{\tilde{C}_2}{2} Z^2 \right) \right] = - \int_0^\infty \exp \left( \frac{\tilde{C}_2}{2} t^2 \right) d[1 - F(t)] \\ = 1 + \int_0^\infty \tilde{C}_2 t \exp \left( \frac{\tilde{C}_2}{2} t^2 \right) [1 - F(t)] dt \\ \leq 1 + \tilde{C}_1 \int_0^\infty \tilde{C}_2 t \exp \left( -\frac{\tilde{C}_2}{2} t^2 \right) dt = 1 + \tilde{C}_1.$$

With the Taylor series of the exponential function, we obtain

$$E \left[ \exp \left( \frac{\tilde{C}_2}{2} Z^2 \right) \right] = \sum_{k=0}^\infty \frac{\tilde{C}_2^k E(Z^{2k})}{2^k k!} \geq \frac{\tilde{C}_2^m E(Z^{2m})}{2^m m!}$$

for any nonnegative integer  $m$ . Thus  $E(Z^{2m}) \leq (1 + \tilde{C}_1)(2\tilde{C}_2^{-1})^m m!$ .



### E.6. Lemma 5 and its proof.

LEMMA 5. *Let  $Z$  be a nonnegative random variable satisfying  $P(Z > t) \leq \tilde{C}_1 \exp(-\tilde{C}_2 t)$  for all  $t > 0$  with  $\tilde{C}_1, \tilde{C}_2 > 0$  some constants. Then*

$$E \left[ \exp \left( \frac{\tilde{C}_2}{2} Z \right) \right] \leq 1 + \tilde{C}_1 \quad \text{and} \quad E(Z^m) \leq (1 + \tilde{C}_1)(2\tilde{C}_2^{-1})^m m!$$

for any nonnegative integer  $m$ .

*Proof.* Let  $F(t)$  be the cumulative distribution function of  $Z$ . Then

$$1 - F(t) = P(Z > t) \leq \tilde{C}_1 \exp(-\tilde{C}_2 t)$$

for all  $t > 0$ . It follows from integration by parts that

$$\begin{aligned} E \left[ \exp \left( \frac{\tilde{C}_2}{2} Z \right) \right] &= - \int_0^\infty \exp \left( \frac{\tilde{C}_2}{2} t \right) d[1 - F(t)] \\ &= 1 + \int_0^\infty \exp \left( \frac{\tilde{C}_2}{2} t \right) [1 - F(t)] dt \\ &\leq 1 + \tilde{C}_1 \int_0^\infty \frac{\tilde{C}_2}{2} \exp \left( -\frac{\tilde{C}_2}{2} t \right) dt = 1 + \tilde{C}_1. \end{aligned}$$

Applying the Taylor series of the exponential function leads to

$$E \left[ \exp \left( \frac{\tilde{C}_2}{2} Z \right) \right] = \sum_{k=0}^{\infty} \frac{\tilde{C}_2^k E(Z^k)}{2^k k!} \geq \frac{\tilde{C}_2^m E(Z^m)}{2^m m!}$$

for any nonnegative integer  $m$ . Thus  $E(Z^m) \leq (1 + \tilde{C}_1)(2\tilde{C}_2^{-1})^m m!$ .

### E.7. Lemma 6 and its proof.

LEMMA 6. *Under Condition 2, both  $\text{dcov}^2(X_k, \mathbf{y})$  and  $\text{dcov}^2(X_k^*, \mathbf{y}^*)$  are uniformly bounded in  $k$ .*

*Proof.* We will show that  $\text{dcov}^2(X_k^*, \mathbf{y}^*)$  are uniformly bounded in  $1 \leq k \leq p$ . Similar arguments apply to prove that  $\text{dcov}^2(X_k, \mathbf{y})$  are also uniformly bounded in  $k$ . Recall that  $\text{dcov}^2(X_k^*, \mathbf{y}^*) = T_{k1} + T_{k2} - 2T_{k3}$  where  $T_{k1} = E[\phi(X_{1k}^*, X_{2k}^*)\psi(\mathbf{y}_1^*, \mathbf{y}_2^*)]$ ,  $T_{k2} = E[\phi(X_{1k}^*, X_{2k}^*)] E[\psi(\mathbf{y}_1^*, \mathbf{y}_2^*)]$ , and  $T_{k3} =$

$E[\phi(X_{1k}^*, X_{2k}^*)\psi(\mathbf{y}_1^*, \mathbf{y}_3^*)]$ . Here  $\phi(X_{1k}^*, X_{2k}^*) = |X_{1k}^* - X_{2k}^*|$  and  $\psi(\mathbf{y}_1^*, \mathbf{y}_2^*) = \|\mathbf{y}_1^* - \mathbf{y}_2^*\|$ . Thus an application of the triangle inequality gives

$$(A.42) \quad 0 \leq \text{dcov}^2(X_k^*, \mathbf{y}^*) \leq |T_{k1}| + |T_{k2}| + 2|T_{k3}|.$$

To prove that  $\text{dcov}^2(X_k^*, \mathbf{y}^*)$  are uniformly bounded in  $k$ , it suffices to show that each term on the right hand side above is uniformly bounded in view of (A.42). As shown in Step 1 of Part 1 in the proof of Theorem 1, under Condition 2 the first quantity  $T_{k1}$  is uniformly bounded in  $1 \leq k \leq p$ . Using similar arguments, we can show that  $T_{k2}$  and  $T_{k3}$  are also uniformly bounded in  $k$ , which completes the proof.

### E.8. Lemma 7 and its proof.

LEMMA 7. *Assume that Conditions 4–5 hold and  $\log p = o(n^{1/2-2\xi})$ . Then with probability at least  $1 - O\{\exp(-\tilde{C}_1 n^{1/2-2\xi})\}$  for some constant  $\tilde{C}_1 > 0$ , it holds that*

$$\min_{|J| \leq s, \mathbf{\Delta} \in \mathbb{R}^{\tilde{p} \times q} \setminus \{\mathbf{0}\}, \|\mathbf{\Delta}_{J^c}\|_{2,1} \leq 3\|\mathbf{\Delta}_J\|_{2,1}} \frac{\|\tilde{\mathbf{X}}\mathbf{\Delta}\|_F}{\sqrt{n}\|\mathbf{\Delta}_J\|_F} \geq \frac{\kappa}{2}.$$

*Proof.* The main idea of the proof is to first introduce an event with a high probability and then derive the desired inequality conditional on that event. Define an event

$$\mathcal{E} = \{\|n^{-1}\tilde{\mathbf{X}}^T\tilde{\mathbf{X}} - \mathbf{\Sigma}\|_\infty < \epsilon\},$$

where  $\|\cdot\|_\infty$  denotes the entrywise matrix  $L_\infty$ -norm and  $0 < \epsilon < 1$  will be specified later. In view of the first part of Condition 4, it follows from Lemmas 4 and 10 that  $P(\mathcal{E}) \geq 1 - \tilde{C}_2 \tilde{p}^2 \exp(-\tilde{C}_3 n^{1/2} \epsilon^2)$  for some constants  $\tilde{C}_2, \tilde{C}_3 > 0$ .

From now on, we condition on the event  $\mathcal{E}$ . By the definition of the Frobenius norm, we have

$$(A.43) \quad n^{-1}\|\tilde{\mathbf{X}}\mathbf{\Delta}\|_F^2 = \text{tr}[\mathbf{\Delta}^T(n^{-1}\tilde{\mathbf{X}}^T\tilde{\mathbf{X}} - \mathbf{\Sigma})\mathbf{\Delta}] + \text{tr}(\mathbf{\Delta}^T\mathbf{\Sigma}\mathbf{\Delta}).$$

For any matrix  $\mathbf{M}$ , denote by  $\mathbf{M}_{ij}$  the  $(i, j)$ -entry of  $\mathbf{M}$ . Then conditional on the event  $\mathcal{E}$ , the first term on the right hand of (A.43) can be bounded

as

$$\begin{aligned}
\text{(A.44)} \quad & \left| \text{tr}[\mathbf{\Delta}^T (n^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} - \mathbf{\Sigma}) \mathbf{\Delta}] \right| = \left| \text{tr}[(n^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} - \mathbf{\Sigma}) \mathbf{\Delta} \mathbf{\Delta}^T] \right| \\
& = \left| \sum_{i=1}^{\tilde{p}} \sum_{j=1}^{\tilde{p}} (n^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} - \mathbf{\Sigma})_{ij} (\mathbf{\Delta} \mathbf{\Delta}^T)_{ij} \right| \\
& \leq \epsilon \sum_{i=1}^{\tilde{p}} \sum_{j=1}^{\tilde{p}} \sum_{k=1}^q |\mathbf{\Delta}_{ik}| |\mathbf{\Delta}_{jk}| = \epsilon \sum_{k=1}^q \left( \sum_{j \in J} |\mathbf{\Delta}_{jk}| + \sum_{j \in J^c} |\mathbf{\Delta}_{jk}| \right)^2 \\
& \leq 2\epsilon \sum_{k=1}^q \left( \sum_{j \in J} |\mathbf{\Delta}_{jk}| \right)^2 + 2\epsilon \sum_{k=1}^q \left( \sum_{j \in J^c} |\mathbf{\Delta}_{jk}| \right)^2,
\end{aligned}$$

where we have used the fact that  $(a+b)^2 \leq 2(a^2+b^2)$  in the last inequality.

By the Cauchy-Schwarz inequality, for any set  $J$  satisfying  $|J| \leq s$  we have

$$\text{(A.45)} \quad \sum_{k=1}^q \left( \sum_{j \in J} |\mathbf{\Delta}_{jk}| \right)^2 \leq \sum_{k=1}^q |J| \sum_{j \in J} \mathbf{\Delta}_{jk}^2 = |J| \cdot \|\mathbf{\Delta}_J\|_F^2 \leq s \|\mathbf{\Delta}_J\|_F^2.$$

For any  $\mathbf{\Delta} \in \mathbb{R}^{\tilde{p} \times q} \setminus \{\mathbf{0}\}$  satisfying  $\|\mathbf{\Delta}_{J^c}\|_{2,1} \leq 3\|\mathbf{\Delta}_J\|_{2,1}$  with  $|J| \leq s$ , similar arguments apply to show that

$$\begin{aligned}
& \sum_{k=1}^q \left( \sum_{j \in J^c} |\mathbf{\Delta}_{jk}| \right)^2 = \sum_{j \in J^c} \sum_{j' \in J^c} \sum_{k=1}^q |\mathbf{\Delta}_{jk}| |\mathbf{\Delta}_{j'k}| \leq \sum_{j \in J^c} \sum_{j' \in J^c} \left( \sum_{k=1}^q \mathbf{\Delta}_{jk}^2 \right)^{1/2} \\
& \quad \cdot \left( \sum_{k=1}^q \mathbf{\Delta}_{j'k}^2 \right)^{1/2} = \sum_{j \in J^c} \sum_{j' \in J^c} \|\mathbf{\Delta}_j\|_2 \cdot \|\mathbf{\Delta}_{j'}\|_2 = \left( \sum_{j \in J^c} \|\mathbf{\Delta}_j\|_2 \right)^2 \\
& = \|\mathbf{\Delta}_{J^c}\|_{2,1}^2 \leq 9\|\mathbf{\Delta}_J\|_{2,1}^2 = 9 \left( \sum_{j \in J} \|\mathbf{\Delta}_j\|_2 \right)^2 \leq 9|J| \left( \sum_{j \in J} \|\mathbf{\Delta}_j\|_2^2 \right) \\
& = 9|J| \|\mathbf{\Delta}_J\|_F^2 \leq 9s \|\mathbf{\Delta}_J\|_F^2,
\end{aligned}$$

which along with (A.44)–(A.45) entails

$$\left| \text{tr}[\mathbf{\Delta}^T (n^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} - \mathbf{\Sigma}) \mathbf{\Delta}] \right| \leq 20s\epsilon \|\mathbf{\Delta}_J\|_F^2.$$

Combining the above inequality with (A.43) and by Condition 5, we obtain

$$\min_{|J| \leq s, \mathbf{\Delta} \in \mathbb{R}^{\tilde{p} \times q} \setminus \{\mathbf{0}\}, \|\mathbf{\Delta}_{J^c}\|_{2,1} \leq 3\|\mathbf{\Delta}_J\|_{2,1}} \left( \frac{\|\tilde{\mathbf{X}}\mathbf{\Delta}\|_F}{\sqrt{n}\|\mathbf{\Delta}_J\|_F} \right)^2 \geq \kappa^2 - 20s\epsilon.$$

It follows from the second part of Condition that  $s \leq \tilde{C}_4 n^\xi$  for some positive constant  $\tilde{C}_4$ . We choose  $\epsilon = 3\kappa^2/(80\tilde{C}_4 n^\xi)$ . Then we have  $\kappa^2 - 20s\epsilon \geq \kappa^2/4$  and for sufficiently large  $n$ ,  $0 < \epsilon < 1$ . Therefore, it holds with probability at least  $1 - O\{\exp(-\tilde{C}_1 n^{1/2-2\xi})\}$  for some constant  $\tilde{C}_1 > 0$  that

$$\min_{|J| \leq s, \mathbf{\Delta} \in \mathbb{R}^{\tilde{p} \times q} \setminus \{\mathbf{0}\}, \|\mathbf{\Delta}_{J^c}\|_{2,1} \leq 3\|\mathbf{\Delta}_J\|_{2,1}} \frac{\|\tilde{\mathbf{X}}\mathbf{\Delta}\|_F}{\sqrt{n}\|\mathbf{\Delta}_J\|_F} \geq \frac{\kappa}{2},$$

which completes the proof.

### E.9. Lemma 8 and its proof.

LEMMA 8. *Assume that Condition 6 and the first part of Condition 4 hold,  $q \leq p$ ,  $\log p = o(n^{1/3})$ , and  $\lambda = c_3 \sqrt{(\log p)/(nq)}$  with  $c_3 > 0$  some large enough constant. Then with probability at least  $1 - O(p^{-c})$  for some positive constant  $c$ , it holds that*

$$\frac{1}{nq} \|\tilde{\mathbf{X}}^T \mathbf{W}\|_{2,\infty} \leq \frac{\lambda}{2}.$$

*Proof.* An application of the union bound leads to

$$(A.46) \quad P(\|\tilde{\mathbf{X}}^T \mathbf{W}\|_{2,\infty} \geq nq\lambda/2) \leq \sum_{j=1}^{\tilde{p}} P \left\{ \sum_{k=1}^q (\tilde{\mathbf{X}}^T \mathbf{W})_{jk}^2 \geq (nq\lambda/2)^2 \right\}$$

for any  $\lambda \geq 0$ , where  $(\tilde{\mathbf{X}}^T \mathbf{W})_{jk}$  is the  $(j, k)$ -entry of  $\tilde{\mathbf{X}}^T \mathbf{W}$ . The key ingredient of the proof is to bound  $P\{\sum_{k=1}^q (\tilde{\mathbf{X}}^T \mathbf{W})_{jk}^2 \geq (nq\lambda/2)^2\}$ . Define  $T_{jk,1} = \sum_{i=1}^n \tilde{\mathbf{X}}_{ij} \mathbf{W}_{ik} I(|\tilde{\mathbf{X}}_{ij}| \leq L)$  and  $T_{jk,2} = \sum_{i=1}^n \tilde{\mathbf{X}}_{ij} \mathbf{W}_{ik} I(|\tilde{\mathbf{X}}_{ij}| > L)$ , where  $L > 0$  will be specified later. Since  $(\tilde{\mathbf{X}}^T \mathbf{W})_{jk} = \sum_{i=1}^n \tilde{\mathbf{X}}_{ij} \mathbf{W}_{ik} =$

$T_{jk,1} + T_{jk,2}$ , we deduce

$$\begin{aligned}
\text{(A.47)} \quad P \left\{ \sum_{k=1}^q (\tilde{\mathbf{X}}^T \mathbf{W})_{jk}^2 \geq (nq\lambda/2)^2 \right\} &\leq P \left\{ \sum_{k=1}^q T_{jk,1}^2 \geq (nq\lambda/4)^2 \right\} \\
&\quad + P \left\{ \sum_{k=1}^q T_{jk,2}^2 \geq (nq\lambda/4)^2 \right\} \\
&\leq \sum_{k=1}^q P \{ |T_{jk,1}| \geq \sqrt{qn}\lambda/4 \} + P \left\{ \sum_{k=1}^q T_{jk,2}^2 \geq (nq\lambda/4)^2 \right\}.
\end{aligned}$$

We will deal with the two terms on the right hand side above separately.

We first bound  $P \{ |T_{jk,1}| \geq \sqrt{qn}\lambda/4 \}$ . By the first part of Condition 4, there exist some positive constants  $a_1$  and  $b_1$  such that

$$P(|\mathbf{v}^T \mathbf{x}_i| > t) \leq a_1 \exp(-b_1 t^2)$$

for any  $\|\mathbf{v}\|_2 = 1$  and  $t > 0$ , where  $\mathbf{x}_i^T = (\mathbf{X}_{i1}, \dots, \mathbf{X}_{ip})$  is the  $i$ th row of the main effect design matrix  $\mathbf{X}$ . Then choosing  $\mathbf{v}$  as a unit vector with the  $j$ th component being 1 gives

$$P(|\mathbf{X}_{ij}| > t) \leq a_1 \exp(-b_1 t^2)$$

for any  $1 \leq i \leq n$ ,  $1 \leq j \leq p$ , and  $t > 0$ . Thus it follows from Lemma 4 that

$$E(\mathbf{X}_{ij}^2) \leq 2(1 + a_1)/b_1 \quad \text{and} \quad E(\mathbf{X}_{ij}^4) \leq 8(1 + a_1)/b_1^2$$

for all  $i$  and  $j$ .

Note that  $\tilde{\mathbf{X}}_{ij} = \mathbf{X}_{ij}$  for  $1 \leq j \leq p$  and  $\tilde{\mathbf{X}}_{ij} = \mathbf{X}_{i\ell} \mathbf{X}_{i\ell'}$  with  $1 \leq \ell < \ell' \leq p$  for  $p+1 \leq j \leq \tilde{p}$ . Thus  $E(\tilde{\mathbf{X}}_{jk}^2)$  are uniformly bounded from above by some positive constant  $\tilde{C}_1$ . Similarly, by Condition 6 and Lemma 5 there exist some positive constants  $a_2$  and  $b_2$  such that  $E(|\mathbf{W}_{ik}|^m) \leq a_2 b_2^m m!$  for any nonnegative integer  $m$  and indices  $i$  and  $k$ . Since  $\tilde{\mathbf{X}}_{ij}$  is independent of  $\mathbf{W}_{ik}$ , we have

$$\begin{aligned}
E \left[ |\tilde{\mathbf{X}}_{ij} \mathbf{W}_{ik} I(|\tilde{\mathbf{X}}_{ij}| \leq L)^m \right] &\leq L^{m-2} E(\tilde{\mathbf{X}}_{ij}^2) E(|\mathbf{W}_{ik}|^m) \\
&\leq m! (L b_2)^{m-2} (2a_2 b_2^2 \tilde{C}_1) / 2
\end{aligned}$$

for each integer  $m \geq 2$ . In view of  $T_{jk,1} = \sum_{i=1}^n \tilde{\mathbf{X}}_{ij} \mathbf{W}_{ik} I(|\tilde{\mathbf{X}}_{ij}| \leq L)$ , applying Bernstein's inequality (Lemma 2.2.11 of [33]) yields

$$\text{(A.48)} \quad P \{ |T_{jk,1}| \geq \sqrt{qn}\lambda/4 \} \leq 2 \exp \left( - \frac{qn\lambda^2}{64a_2 b_2^2 \tilde{C}_1 + 8b_2 L \sqrt{q}\lambda} \right).$$

We next bound  $P\left\{\sum_{k=1}^q T_{jk,2}^2 \geq (nq\lambda/4)^2\right\}$ . By the definition of  $T_{jk,2}$ , it is seen that for each  $j$ , such an event satisfies

$$\left\{\sum_{k=1}^q T_{jk,2}^2 \geq (nq\lambda/4)^2\right\} \subset \left\{|\tilde{\mathbf{X}}_{ij}| > L \text{ for some } 1 \leq i \leq n\right\}.$$

Thus using the union bound, we obtain

$$P\left\{\sum_{k=1}^q T_{jk,2}^2 \geq (nq\lambda/4)^2\right\} \leq \sum_{i=1}^n P\{|\tilde{\mathbf{X}}_{ij}| > L\}.$$

Combining this inequality with (A.46)–(A.48) gives

$$\begin{aligned} \text{(A.49)} \quad P(\|\tilde{\mathbf{X}}^T \mathbf{W}\|_{2,\infty} \geq nq\lambda/2) &\leq 2\tilde{p}q \exp\left(-\frac{qn\lambda^2}{64a_2b_2^2\tilde{C}_1 + 8b_2L\sqrt{q}\lambda}\right) \\ &\quad + \sum_{i=1}^n \sum_{j=1}^p P\{|\mathbf{X}_{ij}| > L\} + \sum_{i=1}^n \sum_{1 \leq \ell < \ell' \leq p} P\{|\mathbf{X}_{i\ell}\mathbf{X}_{i\ell'}| > L\} \\ &\leq qp^2 \exp\left(-\frac{qn\lambda^2}{64a_2b_2^2\tilde{C}_1 + 8b_2L\sqrt{q}\lambda}\right) + a_1np \exp(-b_1L^2) \\ &\quad + a_1np^2 \exp(-b_1L). \end{aligned}$$

Note that  $\lambda = c_3\sqrt{(\log p)/(nq)}$  with some large enough positive constant  $c_3$ . Therefore, setting  $L = \tilde{C}_2\sqrt{n/(\log p)}$  for some large positive constant  $\tilde{C}_2$  ensures that there exists some positive constant  $c_4$  such that

$$\text{(A.50)} \quad P(\|\tilde{\mathbf{X}}^T \mathbf{W}\|_{2,\infty} \geq nq\lambda/2) \leq O(p^{-c}),$$

where we have used the assumption that  $q \leq p$  and  $\log p = o(n^{1/3})$ . This concludes the proof.

### E.10. Additional lemmas.

**LEMMA 9** (Hoeffding's inequality). *Let  $X$  be a real-valued random variable with  $E(X) = 0$ . If  $P(a \leq X \leq b) = 1$  for some  $a, b \in \mathbb{R}$ , then  $E[\exp(tX)] \leq \exp[t^2(b-a)^2/8]$  for any  $t > 0$ .*

**LEMMA 10** (Lemma B.4 of Hao and Zhang [16]). *Let  $Z_1, \dots, Z_n$  be independent random variables with zero mean and  $E[\exp(T_0|Z_i|^\alpha)] \leq A_0$  for*

constants  $T_0, A_0 > 0$  and  $0 < \alpha \leq 1$ . Then there exist some constants  $\tilde{C}_3, \tilde{C}_4 > 0$  such that

$$P\left(\left|n^{-1} \sum_{i=1}^n Z_i\right| > \epsilon\right) \leq \tilde{C}_3 \exp(-\tilde{C}_4 n^\alpha \epsilon^2)$$

for any  $0 < \epsilon \leq 1$ .

<p>DEPARTMENT OF INFORMATION SYSTEMS AND DECISION SCIENCES          MIHAYLO COLLEGE OF BUSINESS AND ECONOMICS          CALIFORNIA STATE UNIVERSITY AT FULLERTON          FULLERTON, CA 92831          USA          E-MAIL: <a href="mailto:yinfeiko@usc.edu">yinfeiko@usc.edu</a></p>	<p>DEPARTMENT OF STATISTICS          UNIVERSITY OF CENTRAL FLORIDA          ORLANDO, FL 32816-2370          USA          E-MAIL: <a href="mailto:daoji.li@ucf.edu">daoji.li@ucf.edu</a></p>
---	---

DATA SCIENCES AND OPERATIONS DEPARTMENT  
 MARSHALL SCHOOL OF BUSINESS  
 UNIVERSITY OF SOUTHERN CALIFORNIA  
 LOS ANGELES, CA 90089  
 USA  
 E-MAIL: [fanyingy@marshall.usc.edu](mailto:fanyingy@marshall.usc.edu)  
[jinchilv@marshall.usc.edu](mailto:jinchilv@marshall.usc.edu)