

**SUPPLEMENTARY MATERIAL TO “INNOVATED
SCALABLE EFFICIENT ESTIMATION IN ULTRA-LARGE
GAUSSIAN GRAPHICAL MODELS”**

BY YINGYING FAN AND JINCHI LV

University of Southern California

This Supplementary Material contains the proofs of Theorem 3, Proposition 1, and additional technical details, as well as an extension of ISEE by incorporating the idea of feature screening.

APPENDIX B: ULTRA-LARGE GRAPH SCREENING

B.1. SIS-assisted ISEE. When the scale of the number of nodes p is ultra large, we can exploit the sure independence screening (SIS) in [16] to reduce the computational cost for each scaled Lasso regression. For each node j in the index set A_l with $1 \leq l \leq L$, the SIS ranks the components of the vector

$$(A.1) \quad \mathbf{w} = (w_k)_{k \in A_l^c} = \mathbf{X}_{A_l^c}^T \mathbf{X}_j$$

obtained by componentwise regression and for any given $\zeta \in (0, 1)$, defines a submodel

$$(A.2) \quad \mathcal{M}_{j,l,\zeta} = \{k \in A_l^c : |w_k| \text{ is among the first } [\zeta n] \text{ largest of all}\},$$

where $[\zeta n]$ denotes the integer part of ζn . Here for simplicity, each node random variable X_j is assumed to have standard deviation one as in [16].

Following [16], based on the reduced model $\mathcal{M}_{j,l,\zeta}$ obtained by the SIS one can construct the SIS-SLasso estimator $\hat{\beta}_{j,l}^*$, which is the scaled Lasso estimator $\hat{\beta}_{j,l}$ as defined in (12) with zero components outside the index set $\mathcal{M}_{j,l,\zeta}$ for β . Similarly as in (16), we define the initial ISEE estimator $\hat{\Omega}_{\text{ISEE,ini}}^*$ as the sample covariance matrix

$$(A.3) \quad \hat{\Omega}_{\text{ISEE,ini}}^* = n^{-1}(\hat{\mathbf{X}}^*)^T \hat{\mathbf{X}}^*,$$

where the estimator $\hat{\mathbf{X}}^*$ for the oracle empirical matrix $\tilde{\mathbf{X}}$ is constructed as in (15) using the SIS-SLasso estimator $\hat{\beta}_{j,l}^*$. Then we can construct the ISEE estimator for the graph $\hat{\Omega}_{\text{ISEE},g}$ and the ISEE estimator with refinement

$\widehat{\Omega}_{\text{ISEE}}$ based on the SIS-assisted initial ISEE estimator $\widehat{\Omega}_{\text{ISEE,ini}}^*$ in (A.3) as described in Section 2.3. Similarly the iterative SIS (ISIS) in [16] can also be applied to improve over the SIS in ultra-large scale problems.

B.2. Technical conditions.

CONDITION 3. *It holds that $p > n$ and $\log p = O(n^\gamma)$ for some constant $0 < \gamma < 1 - 2\kappa$ with κ defined in Condition 4.*

CONDITION 4. *There exist some constants $0 \leq \kappa < 1/2$ and $c_1, c_2, c_3 > 0$ such that for each $j \in A_l$ with $1 \leq l \leq L$, the support of the regression coefficient vector $\beta_{j,l} = (\beta_{jlk})_{k \in A_l^c}$ in (11) admits a decomposition $\text{supp}(\beta_{j,l}) = S_{j,l0} \cup S_{j,l1}$, where for each $k \in S_{j,l0}$, $|\beta_{jlk}| \geq c_1 n^{-\kappa}$ and $|\text{cov}(\beta_{jlk}^{-1} X_j, X_k)| \geq c_2$, and for each $k \in A_l^c$, $|\text{cov}(\sum_{m \in S_{j,l1}} \beta_{jlm} X_m, X_k)| \leq c_3 \lambda$. Moreover, it holds that*

$$(A.4) \quad \max_{j \in A_l, 1 \leq l \leq L} \max \left\{ \sum_{m \in S_{j,l1}} |\beta_{jlm}|, \lambda^{-1} \sum_{m \in S_{j,l1}} \beta_{jlm}^2 \right\} = O(K\lambda).$$

Conditions 3 and 4 are additional assumptions that facilitate the analysis of the SIS-assisted ISEE approach and ensure the sure screening property of the SIS procedure as in [16]. In particular, Condition 3 allows the dimensionality p to increase exponentially with sample size n . Condition 4 is imposed to ensure that the SIS-assisted ISEE estimate can enjoy nice asymptotic properties.

B.3. Theoretical properties. As introduced in Section B.1, to reduce the computational cost we can apply ISEE along with SIS or ISIS in the initial step for ultra-large graph screening. The computational cost can be further reduced if we also apply SIS or ISIS in the refinement step of estimating the link strength. In the refinement step, for each identified link (j, k) we can fit model (11) instead on the union of the supports of the j th and k th rows of $\widehat{\Omega}_{\text{ISEE},g}$, with nodes j and k excluded; see (60) in the proof of Theorem 2 for more details.

The following two theorems characterize the performance of the SIS-assisted ISEE estimators in both the initial step and the refinement step.

THEOREM 4. *Assume that the conditions of Theorem 1 and Conditions 3–4 hold and ζ in (A.2) is at least of order $n^{-\gamma_0}$ with some constant $0 < \gamma_0 < 1 - 2\kappa$. Then the SIS-assisted initial ISEE estimator $\widehat{\Omega}_{\text{ISEE,ini}}^*$ in (A.3) satisfies the same properties as in Theorem 1.*

THEOREM 5. *Under the conditions of Theorems 2 and 4, the ISEE estimators $\widehat{\boldsymbol{\Omega}}_{\text{ISEE},g}$ and $\widehat{\boldsymbol{\Omega}}_{\text{ISEE}}$ based on $\widehat{\boldsymbol{\Omega}}_{\text{ISEE},ini}^*$ in (A.3) satisfy the same properties as in Theorem 2.*

APPENDIX C: PROOFS OF ADDITIONAL MAIN RESULTS

C.1. Proof of Theorem 3. By (16), (35), and the definition of the bias corrected initial ISEE estimator $\widehat{\boldsymbol{\Omega}}_{\text{ISEE},cini}$ in (23) and (24), it suffices to consider the off-block-diagonal entries of the initial ISEE estimator $\widehat{\boldsymbol{\Omega}}_{\text{ISEE},ini}$, that is, the submatrices $(\widehat{\boldsymbol{\Omega}}_{\text{ISEE},ini})_{A_l, A_m}$ with $1 \leq l \neq m \leq L$. The bias of the initial ISEE estimator $\widehat{\boldsymbol{\Omega}}_{\text{ISEE},ini}$ comes from these entries. Note that for each $l \neq m$, $(\widehat{\boldsymbol{\Omega}}_{\text{ISEE},ini})_{A_l, A_m}$ admits the representation in (38). By (54) and (55), we see that the aforementioned bias is incurred by the second and third terms $\boldsymbol{\eta}_2$ and $\boldsymbol{\eta}_3$ in (38).

Due to the symmetry, we focus only on the term $\boldsymbol{\eta}_2$. Examining Part 1 of the proof of Theorem 1, we see that the bias in the term $\boldsymbol{\eta}_2$ is caused only by the additive component

$$(A.5) \quad \widetilde{\mathbf{F}}_2 = -\boldsymbol{\Omega}_{A_l}^T \mathbf{F}_2 \widehat{\boldsymbol{\Omega}}_{A_m},$$

where \mathbf{F}_2 defined in (43) is given by $(n^{-1} \mathbf{X}_{A_l}^T \mathbf{E}_{A_l})^T (\widehat{\mathbf{C}}_{A_m}^{A_l} - \mathbf{C}_{A_m}^{A_l})$, and $\widehat{\mathbf{C}}_{A_m}^{A_l}$ and $\mathbf{C}_{A_m}^{A_l}$ denote submatrices of $\widehat{\mathbf{C}}_{A_m}$ and \mathbf{C}_{A_m} consisting of rows with indices in A_l , respectively. We now add a bias correction term $\widehat{\mathbf{C}}_{A_m}^{A_l} \widehat{\boldsymbol{\Omega}}_{A_m}$ as specified in (24) to $(\widehat{\boldsymbol{\Omega}}_{\text{ISEE},ini})_{A_l, A_m}$, and subsequently to $\widetilde{\mathbf{F}}_2$ given in (A.5). Let us consider the resulting new term

$$(A.6) \quad \widetilde{\mathbf{F}}_2^* = \widetilde{\mathbf{F}}_2 + \widehat{\mathbf{C}}_{A_m}^{A_l} \widehat{\boldsymbol{\Omega}}_{A_m} = \mathbf{F}_4 + \mathbf{F}_5 + \mathbf{C}_{A_m}^{A_l} \boldsymbol{\Omega}_{A_m},$$

where $\mathbf{F}_4 = -[\boldsymbol{\Omega}_{A_l}^T (n^{-1} \mathbf{X}_{A_l}^T \mathbf{E}_{A_l})^T - I_{|A_l|}] (\widehat{\mathbf{C}}_{A_m}^{A_l} - \mathbf{C}_{A_m}^{A_l}) \widehat{\boldsymbol{\Omega}}_{A_m}$ and $\mathbf{F}_5 = \mathbf{C}_{A_m}^{A_l} (\widehat{\boldsymbol{\Omega}}_{A_m} - \boldsymbol{\Omega}_{A_m})$. We study these two terms \mathbf{F}_4 and \mathbf{F}_5 separately.

As in Part 1 of the proof of Theorem 1, we condition on the event $\mathcal{E} \cap (\cap_{3 \leq i \leq 5} \mathcal{E}_i)$ hereafter. Note that $\widehat{\mathbf{C}}_{A_m}^{A_l} - \mathbf{C}_{A_m}^{A_l}$ is exactly the matrix \mathbf{F}_3 introduced therein. In light of the definitions of \mathcal{E} , \mathcal{E}_4 , and \mathcal{E}_5 in (A.40) and (45)–(46), by the facts of $\|\boldsymbol{\Omega}_{A_l}\|_\infty = O(1)$ and $\widehat{\boldsymbol{\Omega}}_{A_m} = O(1)$ it holds uniformly over $1 \leq l \neq m \leq L$ that

$$(A.7) \quad \begin{aligned} \|\mathbf{F}_4\|_\infty &\leq O(1) \|n^{-1} \mathbf{X}_{A_l}^T \mathbf{E}_{A_l} - \boldsymbol{\Omega}_{A_l}^{-1}\|_\infty \|\mathbf{F}_3\|_\infty O(1) \\ &\leq O(\lambda) \cdot O(K^\alpha \lambda) = O(K^\alpha \lambda^2). \end{aligned}$$

Using similar arguments to those in the proof of Lemma 2, we can show that $\|\mathbf{C}_{A_m}^{A_l}\|_\infty = \|\mathbf{-}\boldsymbol{\Omega}_{A_l, A_m} \boldsymbol{\Omega}_{A_m}^{-1}\|_\infty = O(1)$, which along with (A.40) entails

$$(A.8) \quad \|\mathbf{F}_5\|_\infty = O\{\max(K\lambda^2, \lambda)\}.$$

Since $\alpha \leq 1/2$ by Condition 2, it follows from (A.7) and (A.8) that

$$(A.9) \quad \|\mathbf{F}_4 + \mathbf{F}_5\|_\infty \leq O\{\max(K\lambda^2, \lambda)\}.$$

Observe that $\mathbf{C}_{A_m}^{A_l} \boldsymbol{\Omega}_{A_m} = -\boldsymbol{\Omega}_{A_l, A_m} \boldsymbol{\Omega}_{A_m}^{-1} \boldsymbol{\Omega}_{A_m} = -\boldsymbol{\Omega}_{A_l, A_m}$. Therefore, combining (A.6) and (A.9) proves the desired bound for the bias corrected initial ISEE estimator $\widehat{\boldsymbol{\Omega}}_{\text{ISEE,cini}}$ with off-block-diagonal entries

$$(\widehat{\boldsymbol{\Omega}}_{\text{ISEE,cini}})_{A_l, A_m} = -[(\widehat{\boldsymbol{\Omega}}_{\text{ISEE,ini}})_{A_l, A_m} + \widehat{\mathbf{C}}_{A_l}^{A_m} \widehat{\boldsymbol{\Omega}}_{A_l} + \widehat{\mathbf{C}}_{A_m}^{A_l} \widehat{\boldsymbol{\Omega}}_{A_m}];$$

that is, with the same probability bound as in Theorem 1 it holds that

$$\left\| \widehat{\boldsymbol{\Omega}}_{\text{ISEE,cini}} - \boldsymbol{\Omega} \right\|_\infty = O\{\max(K\lambda^2, \lambda)\},$$

which order is in fact $O(\lambda)$ as explained in the proof of Theorem 2.

The second part of Theorem 3, which is graph recovery consistency of the bias corrected initial ISEE estimator $\widehat{\boldsymbol{\Omega}}_{\text{ISEE,cini}}$, can be proved using similar arguments to those in the proof for part a of Theorem 2, by noting that $\widehat{\boldsymbol{\Omega}}_{\text{ISEE},g} = T_\tau(\widehat{\boldsymbol{\Omega}}_{\text{ISEE,cini}})$ and $\omega_0^* = C\lambda$ with $C > 0$ some sufficiently large constant.

C.2. Proof of Theorem 4. We first show that the two events \mathcal{H}_1 and \mathcal{H}_2 defined as

$$(A.10) \quad \mathcal{H}_1 = \bigcap_{j \in A_l, 1 \leq l \leq L} \{S_{j|0} \subset \mathcal{M}_{j|,\zeta}\}$$

and

$$(A.11) \quad \mathcal{H}_2 = \left\{ \max_{j \in A_l, 1 \leq l \leq L} \left\| n^{-1} \mathbf{X}_{A_l^c}^T \sum_{m \in S_{j|1}} \beta_{jlm} \mathbf{X}_m \right\|_\infty \leq O(\lambda) \right\}$$

have large probabilities. The event \mathcal{H}_1 in (A.10) characterizes the sure screening property of the SIS associated with the sets of indices $S_{j|0}$. It is easy to check that Conditions 1–4 in [16] are entailed by our Conditions 1 and 3–4, with \mathcal{M}_* replaced by $S_{j|0}$. In particular, they verified the property C (a concentration property) for Gaussian distributions.

A key observation is that the proof of Theorem 1 in [16] applies equally well to the case where the set of desired effects $S_{j|0}$ plays the role of \mathcal{M}_*

and the set of additional effects $S_{jl1} = \text{supp}(\boldsymbol{\beta}_{j,l}) \setminus S_{jl0}$ may not be empty. Thus an application of the same arguments leads to a similar conclusion to that in Theorem 1 of [16]; that is, for ζ at least in the order of $n^{-\gamma_0}$ with some positive constant $\gamma_0 < 1 - 2\kappa$, we have

$$(A.12) \quad P\{S_{jl0} \subset \mathcal{M}_{j,l,\zeta}\} = 1 - O\left\{\exp[-Cn^{1-2\kappa}/(\log n)]\right\},$$

where C is some positive constant. Since $\log p = O(n^\gamma)$ with constant $0 < \gamma < 1 - 2\kappa$ by Condition 3, we see immediately from (A.12) and Bonferroni's inequality over all nodes j in the index sets A_l that

$$(A.13) \quad P(\mathcal{H}_1) \geq 1 - p \cdot o\left\{p^{-(\delta-1)}\right\} = 1 - o\left\{p^{-(\delta-2)}\right\}.$$

Note that for each $k \in A_l^c$, the expectation of $n^{-1}\mathbf{X}_k^T \sum_{m \in S_{jl1}} \beta_{jlm} \mathbf{X}_m$ is equal to $\text{cov}(\sum_{m \in S_{jl1}} \beta_{jlm} X_m, X_k)$. Thus in view of the assumption of $\max_{k \in A_l^c} |\text{cov}(\sum_{m \in S_{jl1}} \beta_{jlm} X_m, X_k)| \leq c_3\lambda$ by Condition 4, using similar arguments to those for proving (A.35) with t chosen to be $[(\delta+1)(\log p)/(cn)]^{1/2}$ leads to

$$(A.14) \quad P(\mathcal{H}_2) \geq 1 - p(p-1) \cdot O\left\{p^{-(\delta+1)}\right\} = 1 - o\left\{p^{-(\delta-2)}\right\}.$$

Combining (A.13) and (A.14) yields the desired probability bound

$$(A.15) \quad P(\mathcal{H}_1 \cap \mathcal{H}_2) \geq 1 - o\left\{p^{-(\delta-2)}\right\}.$$

From now on we condition on the event $\mathcal{H}_1 \cap \mathcal{H}_2$. On this event, for each node j in the index set A_l , the submodel $\mathcal{M}_{j,l,\zeta}$ given by the SIS contains the set of desired effects S_{jl0} . In light of (A.11), we can treat the component $\sum_{m \in S_{jl1}} \beta_{jlm} \mathbf{X}_m$ of the mean vector $\mathbf{X}_{A_l^c} \boldsymbol{\beta}_{j,l}$ in the univariate linear regression model (11) as part of the error vector in the technical analysis for the scaled Lasso. A key observation is that all the error bounds and probability bounds used in the arguments for proving Lemma 1 hold uniformly over the submodels $\mathcal{M}_{j,l,\zeta}$. Thus an application of the proof of Lemma 1 shows that with probability $1 - o\{p^{-(\delta-2)}\}$ tending to one, it holds uniformly over all nodes j in the index sets A_l with $1 \leq l \leq L$ and all submodels $\mathcal{M}_{j,l,\zeta}$ that

$$(A.16) \quad \left\| \widehat{\boldsymbol{\beta}}_{j,l,\mathcal{M}_{j,l,\zeta}}^* - \boldsymbol{\beta}_{j,l,\mathcal{M}_{j,l,\zeta}} \right\|_1 = O(K\lambda),$$

$$(A.17) \quad n^{-1} \left\| \mathbf{X}_{\mathcal{M}_{j,l,\zeta}} (\widehat{\boldsymbol{\beta}}_{j,l,\mathcal{M}_{j,l,\zeta}}^* - \boldsymbol{\beta}_{j,l,\mathcal{M}_{j,l,\zeta}}) \right\|_2^2 = O(K\lambda^2),$$

where $\widehat{\boldsymbol{\beta}}_{j,l}^*$ denotes the SIS-SLasso estimator, which is the scaled Lasso estimator $\widehat{\boldsymbol{\beta}}_{j,l}$ as defined in (12) with zero components for $\boldsymbol{\beta}$ outside the reduced

index set $\mathcal{M}_{j,l,\zeta}$ obtained by the SIS, and $\mathcal{M}_{j,l,\zeta}$ in the subscripts indicates the corresponding subvectors or submatrices.

In view of (A.15), the intersection of the event $\mathcal{H}_1 \cap \mathcal{H}_2$ and the one given in (A.16)–(A.17) still has large probability $1 - o\{p^{-(\delta-2)}\}$. On such an event, it follows immediately from the sure screening property of $S_{j,l_0} \subset \mathcal{M}_{j,l,\zeta}$, (A.16), and (A.4) that

$$(A.18) \quad \left\| \widehat{\boldsymbol{\beta}}_{j,l}^* - \boldsymbol{\beta}_{j,l} \right\|_1 = O(K\lambda).$$

Note that the proof of Theorem 2 in [33] applies equally well for the largest singular value to show that

$$(A.19) \quad P \left\{ \max_{|\Lambda| \leq \tilde{K}} \lambda_{\max}(n^{-1} \mathbf{X}_{\Lambda}^T \mathbf{X}_{\Lambda}) \leq O(1) \right\} \leq p^{\tilde{K}} e^{-Cn},$$

where \tilde{K} is as defined in the proof of Lemma 1 and C is some positive constant. Since $\tilde{K} \leq \tilde{c}_0 n / (\log p)$ for some sufficiently small positive constant \tilde{c}_0 , it is easy to derive that (A.19) entails

$$(A.20) \quad P \left\{ \max_{|\Lambda| \leq \tilde{K}} \lambda_{\max}(n^{-1} \mathbf{X}_{\Lambda}^T \mathbf{X}_{\Lambda}) \leq O(1) \right\} = 1 - o\{p^{-(\delta-2)}\}.$$

Thus conditioning on this additional event does not change our asymptotic probability bound $1 - o\{p^{-(\delta-2)}\}$.

Denote by $\Lambda_0 = \text{supp}(\boldsymbol{\beta}_{j,l}) \setminus \mathcal{M}_{j,l,\zeta}$. Since $\|\boldsymbol{\beta}_{j,l}\|_0 \leq \tilde{K}$ as shown in the proof of Lemma 1 which implies $|\Lambda_0| \leq \tilde{K}$, by (A.20), (A.4) in Condition 4, and $S_{j,l_0} \subset \mathcal{M}_{j,l,\zeta}$ we have

$$(A.21) \quad \begin{aligned} n^{-1} \|\mathbf{X}_{\Lambda_0} \boldsymbol{\beta}_{j,l,\Lambda_0}\|_2^2 &\leq \lambda_{\max}(n^{-1} \mathbf{X}_{\Lambda_0}^T \mathbf{X}_{\Lambda_0}) \|\boldsymbol{\beta}_{j,l,\Lambda_0}\|_2^2 \\ &\leq \lambda_{\max}(n^{-1} \mathbf{X}_{\Lambda_0}^T \mathbf{X}_{\Lambda_0}) \|\boldsymbol{\beta}_{j,l,S_{j,l_1}}\|_2^2 \\ &\leq O(1) \cdot O(K\lambda^2) = O(K\lambda^2). \end{aligned}$$

Combining (A.17) and (A.21) leads to

$$(A.22) \quad n^{-1} \left\| \mathbf{X}_{A_l^c} (\widehat{\boldsymbol{\beta}}_{j,l}^* - \boldsymbol{\beta}_{j,l}) \right\|_2^2 = O(K\lambda^2).$$

In light of (A.18) and (A.22), we have shown that with probability $1 - o\{p^{-(\delta-2)}\}$ tending to one, it holds uniformly over all nodes j in the index sets A_l with $1 \leq l \leq L$ that the same bounds as (A.23)–(A.24) in Lemma 1 are also valid for the SIS-SLasso estimator. Therefore, the same arguments as in the proof of Theorem 1 carry through.

C.3. Proof of Theorem 5. Theorem 5 holds immediately as a consequence of Theorems 2 and 4.

APPENDIX D: PROOFS OF TECHNICAL RESULTS

D.1. Lemma 1 and its proof.

LEMMA 1. *Under Condition 1, with probability $1 - o\{p^{-(\delta-2)}\}$ tending to one it holds uniformly over all nodes j in the index sets A_l with $1 \leq l \leq L$ and simultaneously that*

$$(A.23) \quad \|\widehat{\boldsymbol{\beta}}_{j,l} - \boldsymbol{\beta}_{j,l}\|_1 = O(K\lambda),$$

$$(A.24) \quad n^{-1} \|\mathbf{X}_{A_l^c}(\widehat{\boldsymbol{\beta}}_{j,l} - \boldsymbol{\beta}_{j,l})\|_2^2 = O(K\lambda^2),$$

$$(A.25) \quad \|n^{-1} \mathbf{X}_{A_l^c}^T \mathbf{E}_{j,l}\|_\infty = O(\lambda),$$

where $\widehat{\theta}_{j,l} = n^{-1} \widehat{\mathbf{E}}_{j,l}^T \widehat{\mathbf{E}}_{j,l}$, $\widetilde{\theta}_{j,l} = n^{-1} \mathbf{E}_{j,l}^T \mathbf{E}_{j,l}$, and the additional subscript l indicates the same scalars and vectors as defined previously with the index set A replaced by A_l .

Proof of Lemma 1. Let us first make a few observations. First, for each index set A_l , the random error vector $\boldsymbol{\eta}_{A_l}$ in the scalar form of the multivariate linear regression model (7) with index set $A = A_l$ is Gaussian with mean $\mathbf{0}$ and covariance matrix $\boldsymbol{\Omega}_{A_l}^{-1}$ and independent of $\mathbf{x}_{A_l^c}$. Since by Condition 1, the spectrum of the precision matrix $\boldsymbol{\Omega}$ is bounded between M^{-1} and M . We see immediately that the spectrum of its principal submatrix $\boldsymbol{\Omega}_{A_l}$ is also bounded between M^{-1} and M , so is that of its inverse $\boldsymbol{\Omega}_{A_l}^{-1}$. This shows that for each corresponding univariate linear regression model (11), its error vector $\mathbf{E}_{j,l}$ is $N(\mathbf{0}, \theta_{j,l} I_n)$ with marginal variance $\theta_{j,l}$ bounded between M^{-1} and M , where the additional subscript l indicates the same scalars and vectors as defined previously with the index set A replaced by A_l .

Second, by Condition 1, the precision matrix $\boldsymbol{\Omega}$ is K -sparse, that is, each of its row or column has at most K nonzero off-diagonal entries. Since $\max_l |A_l| = O(1)$, it follows that the total number of nonzero entries \widetilde{K} in the submatrix $\boldsymbol{\Omega}_{A_l^c, A_l}$ is bounded from above by $K|A_l| = O(K)$. In view of $K \leq c_0 n / (\log p)$ for some sufficiently small positive constant c_0 , we have $\widetilde{K} \leq \widetilde{c}_0 n / (\log p)$ with $\widetilde{c}_0 = O(c_0)$ still some sufficiently small positive constant. Thus for each index set A_l , the regression coefficient matrix $\mathbf{C}_{A_l} = -\boldsymbol{\Omega}_{A_l^c, A_l} \boldsymbol{\Omega}_{A_l}^{-1}$ in the matrix form of the multivariate linear regression model (9) with index set $A = A_l$ satisfies that each column vector has at

most \tilde{K} nonzero components. This shows that for each corresponding univariate linear regression model (11), its regression coefficient vector $\beta_{j,l}$ has sparsity $\|\beta_{j,l}\|_0 \leq \tilde{K} = O(K) \leq \tilde{c}_0 n / (\log p)$ uniformly over all nodes j and index sets A_l .

Third, for each index set A_l , the corresponding univariate linear regression model (11) is a linear regression model with Gaussian design matrix $\mathbf{X}_{A_l^c}$ and Gaussian error vector $\mathbf{E}_{j,l}$ that is independent of $\mathbf{X}_{A_l^c}$. Note that in light of $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$ and (1), $\mathbf{X}_{A_l^c} \sim N(\mathbf{0}, I_n \otimes \Sigma_{A_l^c})$, where $\Sigma_{A_l^c}$ denotes the principal submatrix of Σ given by the index set A_l^c . Since Ω has spectrum bounded between M^{-1} and M , the spectrum of $\Sigma = \Omega^{-1}$ is also bounded between M^{-1} and M and so is that of its principal submatrix $\Sigma_{A_l^c}$.

Denote by $\mathcal{E}_{j,l}$ the event that the bounds (A.23)–(A.25) hold simultaneously for node j in the index set A_l . With the above three observations, an application of the proof of Lemma 2 in [39] shows that

$$(A.26) \quad P(\mathcal{E}_{j,l}) = 1 - o\left\{p^{-(\delta-1)}\right\}.$$

Thus applying Bonferroni's inequality over all nodes j in the index sets A_l along with (A.26) yields the uniform bounds (A.23)–(A.25) satisfied with probability

$$(A.27) \quad P(\mathcal{E}_1) \geq 1 - p \cdot o\left\{p^{-(\delta-1)}\right\} = 1 - o\left\{p^{-(\delta-2)}\right\}$$

which converges to one since $\delta \geq 2$, where the event \mathcal{E}_1 is defined as

$$(A.28) \quad \mathcal{E}_1 = \bigcap_{j \in A_l, 1 \leq l \leq L} \mathcal{E}_{j,l}.$$

In view of $\hat{\mathbf{E}}_{j,l} = \mathbf{X}_j - \mathbf{X}_{A_l^c} \hat{\beta}_{j,l}$, the fact that $\hat{\theta}_{j,l} = n^{-1} \hat{\mathbf{E}}_{j,l}^T \hat{\mathbf{E}}_{j,l}$ follows easily from the definition of the minimizer $(\hat{\beta}_{j,l}, \hat{\theta}_{j,l}^{1/2})$ of the scaled Lasso problem (12).

D.2. Lemma 2 and its proof.

LEMMA 2. *Under Condition 1, with probability $1 - o\{p^{-(\delta-2)}\}$ tending to one it holds uniformly over $1 \leq l \leq L$ that*

$$(A.29) \quad \|\hat{\Omega}_{A_l} - \Omega_{A_l}\|_\infty = O\left\{\max(K\lambda^2, \lambda)\right\},$$

where $\|\cdot\|_\infty$ denotes the entrywise L_∞ -norm of a given matrix.

Proof of Lemma 2. Note that by (13) and (9), we have the following decomposition of the residual matrix

$$(A.30) \quad \widehat{\mathbf{E}}_{A_l} = \mathbf{X}_{A_l} - \mathbf{X}_{A_l^c} \widehat{\mathbf{C}}_{A_l} = \mathbf{E}_{A_l} - \mathbf{X}_{A_l^c} (\widehat{\mathbf{C}}_{A_l} - \mathbf{C}_{A_l}),$$

where $\widehat{\mathbf{C}}_{A_l} = (\widehat{\beta}_{j,l})_{j \in A_l}$ is a $(p - |A_l|) \times |A_l|$ matrix of estimated regression coefficients. Combining (14) and (A.30) yields

$$(A.31) \quad \widehat{\boldsymbol{\Omega}}_{A_l}^{-1} - \boldsymbol{\Omega}_{A_l}^{-1} = n^{-1} \widehat{\mathbf{E}}_{A_l}^T \widehat{\mathbf{E}}_{A_l} - \boldsymbol{\Omega}_{A_l}^{-1} = \boldsymbol{\xi}_1 + \boldsymbol{\xi}_2 + \boldsymbol{\xi}_3,$$

where $\boldsymbol{\xi}_1 = n^{-1} \mathbf{E}_{A_l}^T \mathbf{E}_{A_l} - \boldsymbol{\Omega}_{A_l}^{-1}$, $\boldsymbol{\xi}_2 = -2n^{-1} \mathbf{E}_{A_l}^T \mathbf{X}_{A_l^c} (\widehat{\mathbf{C}}_{A_l} - \mathbf{C}_{A_l})$, and $\boldsymbol{\xi}_3 = n^{-1} (\widehat{\mathbf{C}}_{A_l} - \mathbf{C}_{A_l})^T \mathbf{X}_{A_l^c}^T \mathbf{X}_{A_l^c} (\widehat{\mathbf{C}}_{A_l} - \mathbf{C}_{A_l})$. Let us first consider the last two terms $\boldsymbol{\xi}_2$ and $\boldsymbol{\xi}_3$ conditional on the event \mathcal{E}_1 defined in (A.28). On the event \mathcal{E}_1 , bounds (A.25) and (A.23) control the maximum rowwise L_∞ -norm of matrix $n^{-1} \mathbf{E}_{A_l}^T \mathbf{X}_{A_l^c}$ and maximum columnwise L_1 -norm of matrix $\widehat{\mathbf{C}}_{A_l} - \mathbf{C}_{A_l}$, respectively, which lead to

$$(A.32) \quad \|\boldsymbol{\xi}_2\|_\infty = O(K\lambda^2),$$

where $\|\cdot\|_\infty$ denotes the entrywise L_∞ -norm of a given matrix. An application of the Cauchy-Schwarz inequality along with bound (A.24) results in

$$(A.33) \quad \|\boldsymbol{\xi}_3\|_\infty = O(K\lambda^2).$$

Note that bounds (A.32) and (A.33) are uniform over $1 \leq l \leq L$. It remains to consider the first term $\boldsymbol{\xi}_1$.

As mentioned in the proof of Lemma 1, the spectrum of $\boldsymbol{\Omega}_{A_l}^{-1}$ is bounded between M^{-1} and M . In view of (9) and (7), $n^{-1} \mathbf{E}_{A_l}^T \mathbf{E}_{A_l}$ is the oracle sample covariance matrix estimator for $\boldsymbol{\Omega}_{A_l}^{-1}$. Thus the concentration bounds in [41] and [2], together with Bonferroni's inequality and $\max_l |A_l| = O(1)$, yield for any $t \leq \alpha$,

$$(A.34) \quad P \{ \|\boldsymbol{\xi}_1\|_\infty \leq t \} = 1 - O(e^{-cnt^2}),$$

where c and α are some positive constants. Taking $t = [\delta(\log p)/(cn)]^{1/2}$ in (A.34) and applying Bonferroni's inequality over $1 \leq l \leq L$ lead to

$$(A.35) \quad P(\mathcal{E}_2) \geq 1 - p \cdot O(e^{-cnt^2}) = 1 - O \left\{ p^{-(\delta-1)} \right\} = 1 - o \left\{ p^{-(\delta-2)} \right\},$$

where the event \mathcal{E}_2 is defined as

$$(A.36) \quad \mathcal{E}_2 = \left\{ \max_{1 \leq l \leq L} \left\| n^{-1} \mathbf{E}_{A_l}^T \mathbf{E}_{A_l} - \boldsymbol{\Omega}_{A_l}^{-1} \right\|_\infty \leq t = O(\lambda) \right\}.$$

Therefore, combining (A.31)–(A.34) and (A.35) leads to (A.37)

$$P \left\{ \max_{1 \leq l \leq L} \left\| \widehat{\Omega}_{A_l}^{-1} - \Omega_{A_l}^{-1} \right\|_{\infty} = O \left\{ \max(K\lambda^2, \lambda) \right\} \right\} = 1 - o \left\{ p^{-(\delta-2)} \right\}.$$

We still need to derive the bounds for the matrices $\widehat{\Omega}_{A_l}$.

Let us work with the bound $\left\| \widehat{\Omega}_{A_l}^{-1} - \Omega_{A_l}^{-1} \right\|_{\infty} = O \left\{ \max(K\lambda^2, \lambda) \right\}$. Since $|A_l| = O(1)$, the Frobenius norm $\left\| \widehat{\Omega}_{A_l}^{-1} - \Omega_{A_l}^{-1} \right\|_F = O \left\{ \max(K\lambda^2, \lambda) \right\}$. In light of Condition 1, the quantity $O \left\{ \max(K\lambda^2, \lambda) \right\}$ is bounded above by some sufficiently small positive constant. Then it follows from the matrix perturbation theory (Corollary 6.3.8 of [27]) that

$$\begin{aligned} \lambda_{\min}(\widehat{\Omega}_{A_l}^{-1}) &\geq \lambda_{\min}(\Omega_{A_l}^{-1}) - \left\| \widehat{\Omega}_{A_l}^{-1} - \Omega_{A_l}^{-1} \right\|_F \\ &\geq M^{-1} - O \left\{ \max(K\lambda^2, \lambda) \right\} \geq (2M)^{-1} \end{aligned}$$

for large enough n . The above spectral inequality leads to $\lambda_{\max}(\widehat{\Omega}_{A_l}) = \lambda_{\min}^{-1}(\widehat{\Omega}_{A_l}^{-1}) = O(1)$. Similarly, we can show that $\lambda_{\min}(\widehat{\Omega}_{A_l})$ is also bounded away from zero.

Note a fact that the entrywise L_{∞} -norm of any symmetric positive definite matrix is bounded above by its largest eigenvalue. This claim follows from the facts that each diagonal entry is positive and no larger than the largest eigenvalue and that the 2×2 principal submatrix corresponding to each off-diagonal entry is necessarily nonsingular. Since both Ω_{A_l} and $\widehat{\Omega}_{A_l}$ have spectra bounded away from 0 and ∞ , we see that $\left\| \Omega_{A_l} \right\|_{\infty} = O(1)$ and $\left\| \widehat{\Omega}_{A_l} \right\|_{\infty} = O(1)$, which along with $\max_{1 \leq l \leq L} \left\| \widehat{\Omega}_{A_l}^{-1} - \Omega_{A_l}^{-1} \right\|_{\infty} = O \left\{ \max(K\lambda^2, \lambda) \right\}$ and $\max_l |A_l| = O(1)$ entails

$$(A.38) \quad \left\| \widehat{\Omega}_{A_l} - \Omega_{A_l} \right\|_{\infty} = \left\| \Omega_{A_l} \left(\widehat{\Omega}_{A_l}^{-1} - \Omega_{A_l}^{-1} \right) \widehat{\Omega}_{A_l} \right\|_{\infty} = O \left\{ \max(K\lambda^2, \lambda) \right\}.$$

Therefore, combining (A.27), (A.35), and (A.37)–(A.38) yields

$$(A.39) \quad P(\mathcal{E}) = 1 - o \left\{ p^{-(\delta-2)} \right\},$$

where the event \mathcal{E} is defined as

$$(A.40) \quad \mathcal{E} = \mathcal{E}_1 \cap \mathcal{E}_2 \cap \left\{ \max_{1 \leq l \leq L} \left\| \widehat{\Omega}_{A_l} - \Omega_{A_l} \right\|_{\infty} = O \left\{ \max(K\lambda^2, \lambda) \right\} \right\}.$$

Hereafter we condition on the event \mathcal{E} .

D.3. Proof of Proposition 1. For any $\mathbf{\Omega} \in \mathcal{G}(M, K)$, we know that each row of $\mathbf{\Omega}$ has at most $K + 1$ nonzero components and the spectrum of $\mathbf{\Omega}$ is bounded between M^{-1} and M . Thus it follows easily that for $\mathbf{\Sigma} = \mathbf{\Omega}^{-1}$ and any $\mathbf{u} \neq \mathbf{0}$,

$$(A.41) \quad \|\mathbf{u}\|_\infty = \|\mathbf{\Omega}\mathbf{\Sigma}\mathbf{u}\|_\infty \leq \|\mathbf{\Omega}\|_{\infty, \infty} \|\mathbf{\Sigma}\mathbf{u}\|_\infty,$$

where $\|\cdot\|_{\infty, \infty}$ denotes the operator norm of a matrix induced by the L_∞ -norm. Note that $\|\mathbf{\Omega}\|_{\infty, \infty}$ is the maximum rowwise L_1 -norm of $\mathbf{\Omega}$, which is bounded above by $(K + 1)^{1/2}$ multiplied by the maximum rowwise L_2 -norm of $\mathbf{\Omega}$, thanks to the Cauchy-Schwarz inequality and the fact that each row of $\mathbf{\Omega}$ has L_0 -norm bounded above by $K + 1$. By the definition of the spectral norm, the maximum rowwise L_2 -norm of $\mathbf{\Omega}$ is further bounded above by $\lambda_{\max}(\mathbf{\Omega}) \leq M$, which entails

$$(A.42) \quad \|\mathbf{\Omega}\|_{\infty, \infty} \leq (K + 1)^{1/2} M.$$

Combining (A.41)–(A.42) yields the desired bound $\inf\{\|\mathbf{\Sigma}\mathbf{u}\|_\infty/\|\mathbf{u}\|_\infty : \mathbf{u} \neq \mathbf{0}\} \geq (K + 1)^{-1/2} M^{-1}$.

D.4. Lemma 3 and its proof.

LEMMA 3. *Assume that Conditions 1–2 hold and $K^{1+\alpha}\lambda = o(1)$. Then with probability $1 - o\{p^{-(\delta-2)}\}$ tending to one it holds uniformly over all nodes j in the index sets A_l with $1 \leq l \leq L$ that the L_∞ -norm cone invertibility factor*

$$(A.43) \quad F_{\infty, j, l} = \inf \left\{ \frac{\|\widehat{\mathbf{R}}_{j, l} \mathbf{u}\|_\infty}{\|\mathbf{u}\|_\infty} : \|\mathbf{u}_{S_{j, l}^c}\|_1 \leq \xi \|\mathbf{u}_{S_{j, l}}\|_1 \neq 0 \right\}$$

satisfies $F_{\infty, j, l} \geq c_1 F_\infty$, where $c_1 < 1$ is some positive constant, $S_{j, l}$ denotes the support $\text{supp}(\boldsymbol{\beta}_{j, l})$, and $\widehat{\mathbf{R}}_{j, l} = n^{-1} \mathbf{Y}_{A_l^c}^T \mathbf{Y}_{A_l^c}$ with $\mathbf{Y}_{A_l^c}$ the design matrix $\mathbf{X}_{A_l^c}$ rescaled columnwise to have L_2 -norm $n^{1/2}$ for each column.

Proof of Lemma 3. Let \mathbf{R} be the correlation matrix corresponding to the covariance matrix $\mathbf{\Sigma} = (\sigma_{jk})$. Since the spectrum of $\mathbf{\Sigma}$ is bounded between M^{-1} and M thanks to the same property of $\mathbf{\Omega}$, all diagonal entries σ_{jj} of $\mathbf{\Sigma}$ are also bounded between M^{-1} and M and so are all their reciprocals σ_{jj}^{-1} . Thus the L_1 -norms and L_∞ -norms induced by both linear transformations corresponding to matrices $\mathbf{S} = \text{diag}\{\sigma_{11}^{1/2}, \dots, \sigma_{pp}^{1/2}\}$ and

$\mathbf{S}^{-1} = \text{diag}\{\sigma_{11}^{-1/2}, \dots, \sigma_{pp}^{-1/2}\}$ are equivalent to the original ones. Thus it follows from the identity

$$(A.44) \quad \mathbf{R} = \mathbf{S}^{-1} \mathbf{\Sigma} \mathbf{S}^{-1}$$

that the L_∞ -norm cone invertibility factor F'_∞ with $\mathbf{\Sigma}$ replaced by \mathbf{R} in (20) and the original one F_∞ defined for $\mathbf{\Sigma}$ are within a constant factor of each other. To simplify the notation, we still write F'_∞ as F_∞ which is implicitly understood as the L_∞ -norm cone invertibility factor defined for \mathbf{R} hereafter.

For each node j in the index set A_l , define the population version of the L_∞ -norm cone invertibility factor in (A.43) as

$$(A.45) \quad \tilde{F}_{\infty,j,l} = \inf \left\{ \frac{\|\mathbf{R}_{A_l^c} \mathbf{u}\|_\infty}{\|\mathbf{u}\|_\infty} : \|\mathbf{u}_{S_{j,l}^c}\|_1 \leq \xi \|\mathbf{u}_{S_{j,l}}\|_1 \neq 0 \right\},$$

where $\mathbf{R}_{A_l^c}$ denotes the principal submatrix of \mathbf{R} given by the index set A_l^c . As mentioned in the proof of Lemma 1, $|S_{j,l}| = \|\boldsymbol{\beta}_{j,l}\|_0 \leq \tilde{K} = O(K)$, which together with (20) defined for \mathbf{R} and (A.45) leads to

$$(A.46) \quad \tilde{F}_{\infty,j,l} \geq F_\infty.$$

We will show that the empirical version of the L_∞ -norm cone invertibility factor $F_{\infty,j,l}$ in (A.43) concentrates around its population counterpart $\tilde{F}_{\infty,j,l}$ in (A.45) with overwhelming probability.

Using similar arguments to those for proving (A.35) with t chosen to be $[(\delta + 1)(\log p)/(cn)]^{1/2}$, we can show that

$$(A.47) \quad P(\mathcal{F}) = 1 - o\left\{p^{-(\delta-2)}\right\},$$

where $\mathcal{F} = \{\|\hat{\mathbf{\Sigma}} - \mathbf{\Sigma}\|_\infty \leq O(\lambda)\}$ with $\hat{\mathbf{\Sigma}} = n^{-1} \mathbf{X}^T \mathbf{X}$ and $\|\cdot\|_\infty$ denoting the entrywise L_∞ -norm of a given matrix. Note that $\hat{\mathbf{R}}_{j,l} = n^{-1} \mathbf{Y}_{A_l^c}^T \mathbf{Y}_{A_l}$ is simply the principal submatrix $\hat{\mathbf{R}}_{A_l^c}$ of the sample correlation matrix

$$(A.48) \quad \hat{\mathbf{R}} = \left(\text{diag}\{\hat{\mathbf{\Sigma}}\}\right)^{-1/2} \hat{\mathbf{\Sigma}} \left(\text{diag}\{\hat{\mathbf{\Sigma}}\}\right)^{-1/2}$$

given by the index set A_l^c . By some standard calculations, we can show that on the event \mathcal{F} , it also holds that $\|\hat{\mathbf{R}} - \mathbf{R}\|_\infty \leq O(\lambda)$. This result together with (A.47) yields

$$(A.49) \quad P(\mathcal{F}_1) \geq 1 - o\left\{p^{-(\delta-2)}\right\},$$

where the event \mathcal{F}_1 is defined as the intersection of events \mathcal{F} and $\{\|\widehat{\mathbf{R}} - \mathbf{R}\|_\infty \leq O(\lambda)\}$.

Finally let us do some algebraic calculations conditional on event \mathcal{F}_1 . On this event, for each $\mathbf{u} \in \mathbb{R}^{p-|A_l|}$ satisfying $\|\mathbf{u}_{S_{j,l}^c}\|_1 \leq \xi \|\mathbf{u}_{S_{j,l}}\|_1 \neq 0$ we have

$$\begin{aligned}
(A.50) \quad \|\widehat{\mathbf{R}}_{j,l} \mathbf{u}\|_\infty &= \|\widehat{\mathbf{R}}_{A_l^c} \mathbf{u}\|_\infty \geq \|\mathbf{R}_{A_l^c} \mathbf{u}\|_\infty - \left\| \left(\widehat{\mathbf{R}}_{A_l^c} - \mathbf{R}_{A_l^c} \right) \mathbf{u} \right\|_\infty \\
&\geq \widetilde{F}_{\infty,j,l} \|\mathbf{u}\|_\infty - \|\widehat{\mathbf{R}} - \mathbf{R}\|_\infty \|\mathbf{u}\|_1 \\
&\geq \widetilde{F}_{\infty,j,l} \|\mathbf{u}\|_\infty - O(\lambda)(1 + \xi) \|\mathbf{u}_{S_{j,l}}\|_1 \\
&\geq \widetilde{F}_{\infty,j,l} \|\mathbf{u}\|_\infty - O(\lambda)(1 + \xi) |S_{j,l}| \|\mathbf{u}_{S_{j,l}}\|_\infty \\
&\geq \left[\widetilde{F}_{\infty,j,l} - O(K\lambda) \right] \|\mathbf{u}\|_\infty,
\end{aligned}$$

since $|S_{j,l}| \leq \widetilde{K} = O(K)$. Therefore, combining (A.46), (A.49)–(A.50), and the assumption of $K^{1+\alpha}\lambda = o(1)$ yields $F_{\infty,j,l} \geq c_1 F_\infty$ for some positive constant $c_1 < 1$, uniformly over all nodes j in the index sets A_l with $1 \leq l \leq L$.

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