

# Joint estimation of precision matrices in heterogeneous populations\*

Takumi Saegusa

*Department of Mathematics, University of Maryland*

*College Park, MD 20742 USA*

*e-mail: [tsaegusa@math.umd.edu](mailto:tsaegusa@math.umd.edu)*

and

Ali Shojaie

*Department of Biostatistics, University of Washington*

*Seattle, WA 98195 USA*

*e-mail: [ashojaie@uw.edu](mailto:ashojaie@uw.edu)*

**Abstract:** We introduce a general framework for estimation of inverse covariance, or precision, matrices from heterogeneous populations. The proposed framework uses a Laplacian shrinkage penalty to encourage similarity among estimates from disparate, but related, subpopulations, while allowing for differences among matrices. We propose an efficient alternating direction method of multipliers (ADMM) algorithm for parameter estimation, as well as its extension for faster computation in high dimensions by thresholding the empirical covariance matrix to identify the joint block diagonal structure in the estimated precision matrices. We establish both variable selection and norm consistency of the proposed estimator for distributions with exponential or polynomial tails. Further, to extend the applicability of the method to the settings with unknown populations structure, we propose a Laplacian penalty based on hierarchical clustering, and discuss conditions under which this data-driven choice results in consistent estimation of precision matrices in heterogeneous populations. Extensive numerical studies and applications to gene expression data from subtypes of cancer with distinct clinical outcomes indicate the potential advantages of the proposed method over existing approaches.

**Keywords and phrases:** Graph Laplacian, graphical modeling, heterogeneous populations, hierarchical clustering, high-dimensional estimation, precision matrix, sparsity.

Received January 2015.

## Contents

1	Introduction . . . . .	1342
2	Model and estimator . . . . .	1345
2.1	Problem setup . . . . .	1345
2.2	The Laplacian shrinkage estimator . . . . .	1345

---

\*This work was partially supported by NSF grants DMS-1161565 & DMS-1561814 to AS.

2.3	Connections to other estimators . . . . .	1348
3	Theoretical properties . . . . .	1350
3.1	Consistency in spectral norm . . . . .	1351
3.2	Model selection consistency . . . . .	1352
3.3	Additional results . . . . .	1354
4	Laplacian shrinkage based on hierarchical clustering . . . . .	1355
5	Algorithms . . . . .	1357
6	Numerical results . . . . .	1360
6.1	Simulation experiments . . . . .	1360
6.1.1	Known subpopulation network $G$ . . . . .	1360
6.1.2	Unknown subpopulation network $G$ . . . . .	1361
6.2	Genetic networks of cancer subtypes . . . . .	1362
7	Discussion . . . . .	1363
8	Appendix: Proofs and technical details . . . . .	1365
8.1	Consistency in matrix norms . . . . .	1365
8.2	Model selection consistency . . . . .	1374
	References . . . . .	1389

## 1. Introduction

Estimation of large inverse covariance, or precision, matrices has received considerable attention in recent years. This interest is in part driven by the advent of high-dimensional data in many scientific areas, including high throughput *omics* measurements, functional magnetic resonance images (fMRI), and applications in finance and industry. Applications of various statistical methods in such settings require an estimate of the (inverse) covariance matrix. Examples include dimension reduction using principal component analysis (PCA), classification using linear or quadratic discriminant analysis (LDA/QDA), and discovering conditional independence relations in Gaussian graphical models (GGM).

In high-dimensional settings, where the data dimension  $p$  is often comparable or larger than the sample size  $n$ , regularized estimation procedures often result in more reliable estimates. Of particular interest is the use of sparsity inducing penalties, specifically the  $\ell_1$  or lasso penalty [30], which encourages sparsity in off-diagonal elements of the precision matrix [34, 7, 8, 33]. Theoretical properties of  $\ell_1$ -penalized precision matrix estimation have been studied under both multivariate normality, as well as some relaxations of this assumption [19, 26, 4, 25].

Sparse estimation is particularly relevant in the setting of GGMs, where conditional independencies among variables correspond to zero off-diagonal elements of the precision matrix [14]. The majority of existing approaches for estimation of high-dimensional precision matrices, including those cited in the previous paragraph, assume that the observations are identically distributed, and correspond to a single population. However, data sets in many application areas include observations from several distinct subpopulations. For instance, gene expression measurements are often collected for both healthy subjects, as well as patients diagnosed with different subtypes of cancer. Despite increasing

evidence for differences among genetic networks of cancer and healthy subjects [11, 27], the networks are also expected to share many common edges. Separate estimation of graphical models for each of the subpopulations would ignore the common structure of the precision matrices, and may thus be inefficient; this inefficiency can be particularly significant in high-dimensional low sample settings, where  $p \gg n$ .

To address the need for estimation of graphical models in related subpopulations, few methods have been recently proposed for joint estimation of  $K$  precision matrices  $\Omega^{(k)} = (\omega_{ij}^{(k)})_{i,j=1}^p \in \mathbb{R}^{p \times p}$ ,  $k = 1, \dots, K$  [9, 6]. These methods extend the penalized maximum likelihood approach by combining the Gaussian likelihoods for the  $K$  subpopulations

$$\ell_n(\Omega) = \frac{1}{n} \sum_{k=1}^K n_k \left( \log \det(\Omega^{(k)}) - \text{tr} \left( \hat{\Sigma}_n^{(k)} \Omega^{(k)} \right) \right). \quad (1)$$

Here,  $n_k$  and  $\hat{\Sigma}_n^{(k)}$  are the number of observations and the sample covariance matrix for the  $k$ th subpopulation, respectively,  $n = \sum_{k=1}^K n_k$  is the total sample size and  $\text{tr}(\cdot)$  and  $\det(\cdot)$  denote matrix trace and determinant.

To encourage similarity among estimated precision matrices, Guo et al. [9] modeled the  $(i, j)$ -element of  $\Omega^{(k)}$  as product of a common factor  $\theta_{ij}$  and group-specific parameters  $\gamma_{ij}^{(k)}$ , i.e.  $\omega_{ij}^{(k)} = \delta_{ij} \gamma_{ij}^{(k)}$ . Identifiability of the estimates is ensured by assuming  $\delta_{ij} \geq 0$ . A zero common factor  $\delta_{ij} = 0$  induces sparsity across all subpopulations, whereas  $\gamma_{ij}^{(k)} = 0$  results in condition-specific sparsity for  $\omega_{ij}^{(k)}$ . This reparametrization results in a non-convex optimization problem based on the Gaussian likelihood with  $\ell_1$ -penalties  $\sum_{i \neq j} \delta_{ij}$  and  $\sum_{i \neq j} \sum_{k=1}^K |\gamma_{ij}^{(k)}|$ . Danaher et al. [6] proposed two alternative estimators by adding an additional convex penalty to the graphical lasso objective function: either a fused lasso penalty  $\sum_{i \neq j} \sum_{k \neq k'} |\omega_{ij}^{(k)} - \omega_{ij}^{(k')}|$  (FGL), or a group lasso penalty  $\sum_{i \neq j} \sqrt{\sum_{k=1}^K (\omega_{ij}^{(k)})^2}$  (GGL). The fused lasso penalty has also been used by Kolar et al. [13], for joint estimation of multiple graphical models in multiple time points. The fused lasso penalty strongly encourages the values of  $\omega_{ij}^{(k)}$  to be similar across all subpopulations, both in values as well as sparsity patterns. On the other hand, the group lasso penalty results in similar estimates by shrinking all  $\omega_{ij}^{(k)}$  across subpopulations to zero if  $\sum_{k=1}^K (\omega_{ij}^{(k)})^2$  is small.

Despite their differences, methods of Guo et al. [9] and Danaher et al. [6] inherently assume that precision matrices in  $K$  subpopulations are equally similar to each other, in that they encourage  $\omega_{ij}^{(k)}$  and  $\omega_{ij}^{(k')}$  and  $\omega_{ij}^{(k)}$  and  $\omega_{ij}^{(k'')}$  to be equally similar. However, when  $K > 2$ , some subpopulations are expected to be more similar to each other than others. For instance, it is expected that genetic networks of two subtypes of cancer be more similar to each other than to the network of normal cells. Similarly, differences among genetic networks of various strains of a virus or bacterium are expected to correspond to the evolutionary lineages of their phylogenetic trees. Unfortunately, existing methods for joint estimation of multiple graphical models ignore this heterogeneity in

multiple subpopulations. Furthermore, existing methods assume subpopulation memberships are known, which limits their applicability in settings with complex but *unknown* population structures; an important example is estimation of genetic networks of cancer cells with unknown subtypes.

In this paper, we propose a general framework for joint estimation of multiple precision matrices by capturing the heterogeneity among subpopulations. In this framework, similarities among disparate subpopulations are presented using a *subpopulation network*  $G(V, E, W)$ , a weighted graph whose node set  $V$  is the set of subpopulations. The edges in  $E$  and the weights  $W_{kk'}$  for  $(k, k') \in E$  represent the degree of similarity between any two subpopulations  $k, k'$ . In the special case where  $W_{kk'} = 1$  for all  $k, k'$ , the subpopulation similarities are only captured by the structure of the graph  $G$ . An example of such a subpopulation network is the line graph corresponding to observations over multiple time points, which is used in estimation of time-varying graphical models [13]. As we will show in Section 2.3, other existing methods for joint estimation of multiple graphical models, e.g. proposals of Danaher et al. [6], can also be seen as special cases of this general framework.

Our proposed estimator is the solution to a convex optimization problem based on the Gaussian likelihood with both  $\ell_1$  and graph Laplacian [15] penalties. The graph Laplacian has been used in other applications for incorporating *a priori* knowledge in classification [24], for principal component analysis on network data [28], and for penalized linear regression with correlated covariates [15, 10, 32, 18, 17, 37]. The Laplacian penalty encourages similarity among estimated precision matrices according to the subpopulation network  $G$ . The  $\ell_1$ -penalty, on the other hand, encourages sparsity in the estimated precision matrices. Together, these two penalties capture both unique patterns specific to each subpopulation, as well as common patterns shared among different subpopulations.

We first discuss the setting where  $G(V, E, W)$  is known from external information, e.g. known phylogenetic trees (Section 2), and later discuss the estimation of the subpopulation memberships and similarities using hierarchical clustering (Section 4). We propose an alternating methods of multipliers (ADMM) algorithm [3] for parameter estimation, as well as its extension for efficient computation in high dimensions by decomposing the problem into block-diagonal matrices. Although we use the Gaussian likelihood, our theoretical results also hold for non-Gaussian distributions. We establish model selection and norm consistency of the proposed estimator under different model assumptions (Section 3), with improved rates of convergence over existing methods based on penalized likelihood. We also establish the consistency of the proposed algorithm for the estimation of multiple precision matrices, in settings where the subpopulation network  $G$  or subpopulation memberships are unknown. To achieve this, we establish the consistency of hierarchical clustering in high dimensions, by generalizing recent results of Borysov et al. [1] to the setting of arbitrary covariance matrices, which is of independent interest.

The rest of the paper is organized as follows. In Section 2 we describe the formal setup of the problem and present our estimator. Theoretical properties

of the proposed estimator are studied in Section 3, and Section 4 discusses the extension of the method to the setting where the subpopulation network is unknown. The ADMM algorithm for parameter estimation and its extension for efficient computation in high dimensions are presented in Section 5. Results of the numerical studies, using both simulated and real data examples, are presented in Section 6. Section 7 concludes the paper with a discussion. Technical proofs are collected in the Appendix.

## 2. Model and estimator

### 2.1. Problem setup

Consider  $K$  subpopulations with distributions  $\mathcal{P}^{(k)}$ ,  $k = 1, \dots, K$ . Let  $X^{(k)} = (X^{(k),1}, \dots, X^{(k),p})^T \in \mathbb{R}^p$  be a random vector from the  $k$ th subpopulation with mean  $\mu_k$  and the covariance matrix  $\Sigma_0^{(k)} = (\sigma_{ij}^{(k)})_{i,j=1}^p$ . Suppose that an observation comes from the  $k$ th subpopulation with probability  $\pi_k > 0$ .

Our goal is to estimate the precision matrices  $\Omega_0^{(k)} \equiv (\Sigma_0^{(k)})^{-1} \equiv (\omega_{ij}^{(k)})_{i,j=1}^p$ ,  $k = 1, \dots, K$ . To this end, we use the Gaussian log-likelihood based on the *correlation matrix* (see Rothman et al. [26]) as a working model for estimation of true  $\Omega_0^{(k)}$ ,  $k = 1, \dots, K$ . Let  $X_i^{(k)}$ ,  $i = 1, \dots, n_k$ , be independent and identically distributed (i.i.d.) copies from  $\mathcal{P}^{(k)}$ ,  $k = 1, \dots, K$ . We denote the correlation matrices and their inverse by  $\Theta^{(k)} = (\theta_{ij}^{(k)})_{i,j=1}^p$ , and  $\Psi^{(k)} = (\psi_{ij}^{(k)})_{i,j=1}^p$ ,  $k = 1, \dots, K$ , respectively. The Gaussian log-likelihood based on the correlation matrix can then be written as

$$\tilde{\ell}_n(\Theta) = \frac{1}{n} \sum_{k=1}^K n_k \left( \log \det(\Theta^{(k)}) - \text{tr} \left( \Psi_n^{(k)} \Theta^{(k)} \right) \right), \quad (2)$$

where  $\Psi_n^{(k)}$ ,  $k = 1, \dots, K$  is the sample correlation matrix for subpopulation  $k$ .

Examining the derivative of (2), which consists of  $\Psi_0^{(k)} - \Psi_n^{(k)}$ ,  $k = 1, \dots, K$ , justifies its use as a working model for non-Gaussian data: the stationary points of (2) is  $\Psi_n^{(k)}$ , which gives a consistent estimate of  $\Psi_0^{(k)}$ . Thus we do not, in general, need to assume multivariate normality. However, in certain applications, for instance LDA/QDA and GGM, the resulting estimate is useful only if the data follows a multivariate normal distribution.

### 2.2. The Laplacian shrinkage estimator

Let  $\Theta = (\Theta^{(1)}, \dots, \Theta^{(K)})$  and write  $\Theta_{ij} = (\theta_{ij}^{(1)}, \dots, \theta_{ij}^{(K)})^T \in \mathbb{R}^K$ ,  $i, j = 1, \dots, p$  for a vector of  $(i, j)$ -elements across subpopulations. Our proposed estimator, Laplacian Shrinkage for Inverse Covariance matrices from Heterogeneous populations (LASICH), first estimates the inverse of the correlation matrices for each of the  $K$  subpopulations, and then transforms them into the estimator of inverse covariance matrices, as in Rothman et al. [26]. In particular, we first obtain

the estimate  $\hat{\Theta}$  of the true inverse correlation matrix by solving the following optimization problem

$$\begin{aligned}\hat{\Theta}_{\rho_n} &\equiv \arg \min_{\Theta=\Theta^T, \Theta \succ 0} -\tilde{\ell}_n(\Theta) + \rho_n \|\Theta\|_1 + \rho_n \rho_2 \|\Theta\|_L \\ &\equiv \arg \min_{\Theta=\Theta^T, \Theta \succ 0} -\tilde{\ell}_n(\Theta) + \rho_n \sum_{k=1}^K \sum_{i \neq j} |\Theta_{ij}^{(k)}| + \rho_n \rho_2 \sum_{i \neq j} \|\Theta_{ij}\|_L, \quad (3)\end{aligned}$$

where  $\Theta = \Theta^T$  enforces the symmetry of individual inverse correlation matrices, i.e.  $\Theta^{(k)} = (\Theta^{(k)})^T$ , and  $\Theta \succ 0$  requires that  $\Theta^{(k)}$  is positive definite for  $k = 1, \dots, K$ . The  $\ell_1$ -penalty  $\|\Theta\|_1 = \sum_{k=1}^K \|\Theta^{(k)}\|_1$  in (3) encourages sparsity in estimated inverse correlation matrices. The graph Laplacian penalty, on the other hand, exploits the information in the subpopulation network  $G$  to encourage similarity among values of  $\theta_{ij}^{(k)}$  and  $\theta_{ij}^{(k')}$ . The tuning parameters  $\rho_n$  and  $\rho_n \rho_2$  control the size of each penalty term.

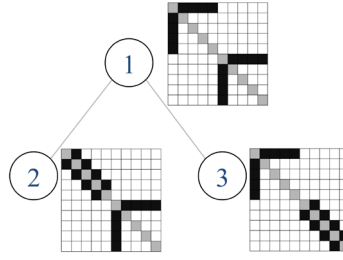


FIG 1. Illustration of similarities in the sparsity patterns of precision matrices  $\Omega^{(1)}, \Omega^{(2)}$  and  $\Omega^{(3)}$ . Nonzero and zero off-diagonal entries are colored in black and white, respectively, while diagonal entries are colored in gray. The associated subpopulation network  $G$  reflects the similarities between precision matrices of subpopulations 1 and 2 and 1 and 3. The simulation experiments in Section 6.1 use a similar subpopulation network in a high-dimensional setting.

Figure 1 illustrates the motivation for the graph Laplacian penalty  $\|\Theta_{ij}\|_L$  in (3). The gray-scale images in the figure show the hypothetical sparsity patterns of precision matrices  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$  for three related subpopulations. Here,  $\Theta^{(1)}$  consists of two blocks with one “hub” node in each block; in  $\Theta^{(2)}$  and  $\Theta^{(3)}$  one of the blocks is changed into a “banded” structure. It can be seen that one of the two blocks in both  $\Theta^{(2)}$  and  $\Theta^{(3)}$  have a similar sparsity pattern as  $\Theta^{(1)}$ . However,  $\Theta^{(2)}$  and  $\Theta^{(3)}$  are not similar. The subpopulation network  $G$  in this figure captures the relationship among precision matrices of the three subpopulations. Such complex relationships cannot be captured using the existing approaches, e.g. Guo et al. [9], Danaher et al. [6], which encourage all precision matrices to be equally similar to each other. More generally,  $G$  can be a weighted graph,  $G(V, E, W)$ , whose nodes represent the subpopulations  $1, \dots, K$ . The edge weights  $W : E \rightarrow \mathbb{R}_+$  represent the similarity among pairs of subpopulations, with larger values of  $W_{kk'} \equiv W(k, k') > 0$  corresponding to more similarity between precision matrices of subpopulations  $k$  and  $k'$ .

In this section, we assume that the weighted graph  $G$  is externally available, and defer the discussion of data-driven choices of  $G$ , based on hierarchical clustering, to Section 4. Given  $G$ , the (unnormalized) graph Laplacian penalty  $\|\Theta_{ij}\|_L$  is defined as

$$\|\Theta_{ij}\|_L = \left\{ \sum_{k,k'=1}^K W_{kk'} \left( \theta_{ij}^{(k)} - \theta_{ij}^{(k')} \right)^2 \right\}^{1/2} \quad (4)$$

where  $W_{kk'} = 0$  if  $k$  and  $k'$  are not connected. The Laplacian shrinkage penalty can be alternatively written as  $\|\Theta_{ij}\|_L = \Theta_{ij}^T L \Theta_{ij}$ , where  $L = (l_{kk'})_{k,k'=1}^K \in \mathbb{R}^{K \times K}$  is the Laplacian matrix [5] of the subpopulation network  $G$  defined as

$$l_{kk'} = \begin{cases} d_k - W_{kk}, & k = k', d_k \neq 0, \\ -W_{kk'}, & k \neq k', \\ 0, & \text{otherwise,} \end{cases}$$

where  $d_k = \sum_{k' \neq k} W_{kk'}$  is the degree of node  $k$  in  $G$  with  $W_{kk'} = 0$  if  $k$  and  $k'$  are not connected. The Laplacian shrinkage penalty can also be defined in terms of the *normalized* graph Laplacian,  $I - D^{-1/2} W D^{-1/2}$ , where  $D = \text{diag}(d_1, \dots, d_K)$  is the diagonal degree matrix. The normalized Laplacian penalty,

$$\|\Theta_{ij}\|_L = \left\{ \sum_{k,k'=1}^K W_{kk'} \left( \frac{\theta_{ij}^{(k)}}{\sqrt{d_k}} - \frac{\theta_{ij}^{(k')}}{\sqrt{d_{k'}}} \right)^2 \right\}^{1/2},$$

which we also denote as  $\|\Theta_{ij}\|_L$ , imposes smaller shrinkage on coefficients associated with highly connected subpopulations. We henceforth primarily focus on the normalized penalty.

Given estimates of the inverse correlation matrices  $\hat{\Theta}^{(1)}, \dots, \hat{\Theta}^{(K)}$  from (3), we obtain estimates of precision matrices  $\Omega^{(k)}$  by noting that  $\Omega^{(k)} = \Xi^{(k)} \Theta^{(k)} \Xi^{(k)}$ , where  $\Xi^{(k)}$  is the diagonal matrix of reciprocals of the standard deviations  $\Xi^{(k)} = \text{diag}(\{\sigma_{11}^{(k)}\}^{-1/2}, \dots, \{\sigma_{pp}^{(k)}\}^{-1/2})$ . Our estimator  $\hat{\Omega}_{\rho_n} = (\hat{\Omega}_{\rho_n}^{(1)}, \dots, \hat{\Omega}_{\rho_n}^{(K)})$  of precision matrices  $\Omega$  is thus defined as

$$\hat{\Omega}_{\rho_n}^{(k)} = \{\hat{\Xi}^{(k)}\}^{-1} \hat{\Theta}_{\rho_n}^{(k)} \{\hat{\Xi}^{(k)}\}^{-1}, \quad k = 1, \dots, K,$$

where  $\hat{\Xi}^{(k)} = \text{diag}(1/\{\hat{\sigma}_{11}^{(k)}\}^{1/2}, \dots, 1/\{\hat{\sigma}_{pp}^{(k)}\}^{1/2})$  with sample variance  $\hat{\sigma}_{ii}^{(k)}$  for the  $i$ th element in the  $k$ th subpopulation.

A number of alternative strategies can be used instead of the graph Laplacian penalty in (3). First, similarity among coefficients of precision matrices can also be imposed using a ridge-type penalty,  $\|\Theta_{ij}\|_L^2$ . The main difference is that our penalty  $\|\Theta_{ij}\|_L$  discourages the inclusion of edges  $\theta_{ij}^{(1)}, \dots, \theta_{ij}^{(K)}$  if they are very different across the  $K$  subpopulations. Another option is to use the graph trend filtering [31], which impose a fused lasso penalty over the subpopulation graph  $G$ . Finally, ignoring the weights  $W_{kk'}$  in (4), the Laplacian shrinkage penalty

resembles the Markov random field (MRF) prior used in Bayesian variable selection with structured covariates [16]. While our paper was under review, we became aware of the recent work by Peterson et al. [23], who utilize an MRF prior to develop a Bayesian framework for estimation of multiple Gaussian graphical models. This method assumes that edges between pairs of random variable are formed independently, and is hence more suited for Erdős-Rényi networks. Our penalized estimation framework can be seen as an alternative to using an MRF prior to estimate the precision matrices in a mixture of Gaussian distributions.

### 2.3. Connections to other estimators

To connect our proposed estimator to existing methods for joint estimation of multiple graphical models, we first give an alternative interpretation of the graph Laplacian penalty  $\|\Theta_{ij}\|_L = (\Theta_{ij}^T L \Theta_{ij})^{1/2}$  as a norm for a transformed version of  $\theta_{ij}^{(k)}$ s. More specifically, consider the mapping  $g_G : \mathbb{R}^K \rightarrow \mathbb{R}^K$  defined based on the Laplacian matrix for graph  $G$

$$g_G(\Theta_{ij}) = \begin{cases} 0, & k = k', \\ \sqrt{W_{kk'}} \left( \frac{\theta_{ij}^{(k)}}{\sqrt{2d_k}} - \frac{\theta_{ij}^{(k')}}{\sqrt{2d_{k'}}} \right), & k \neq k', \end{cases}$$

if  $G$  has at least one edge. For a graph with no edges, define  $g_G(\Theta_{ij}) = I_K \otimes \Theta_{ij} = \text{diag}(\Theta_{ij})$ , where  $I_K$  is the  $K$ -identity matrix, and  $\otimes$  denotes the Kronecker product. It can then be seen that the graph Laplacian penalty can be rewritten as

$$\|\Theta_{ij}\|_L = \|g_G(\Theta_{ij})\|_F.$$

where  $\|\cdot\|_F$  is the Frobenius norm.

Using the above interpretation, other methods for joint estimation of multiple graphical models can be seen as penalties on transformations  $g_G(\Theta_{ij})$  corresponding to different graphs  $G$ . We illustrate this connection using the hypothetical subpopulation network shown in Figure 2a.

Consider first the FGL penalty of Danaher et al. [6], applied to elements of the inverse correlation matrix  $|\theta_{ij}^{(k)} - \theta_{ij}^{(k')}|$ . Let  $G_C$  be a complete unweighted graph ( $W_{kk'} = 1 \forall k \neq k'$ ), in which all  $\binom{K}{2}$  node-pairs are connected to each other (Figure 2b). It is then easy to see that

$$\sum_{k \neq l} |\theta_{ij}^{(k)} - \theta_{ij}^{(l)}| = \sqrt{2(K-1)} \|g_{G_C}(\Theta_{ij})\|_1,$$

where the factor of  $\sqrt{2(K-1)}$  can be absorbed into the tuning parameter for the FGL penalty. A similar argument can also be applied to the GGL penalty of Danaher et al. [6],  $\|\Theta_{ij}\|$ , by considering instead an empty graph  $G_e$  with no edges between nodes (Figure 2c). In this case, the mapping  $g_G$  would give a diagonal matrix with elements  $\theta_{ij}^{(k)}$ , and hence  $\|\Theta_{ij}\| = \|g_{G_e}(\Theta_{ij})\|_F$ .



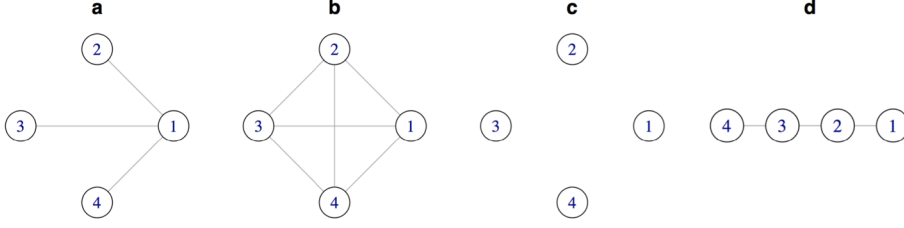


FIG 2. Comparison of subpopulation networks used in the penalty for different methods for joint estimation of multiple precision matrices: **a)** the true network, modeled by LASICH; **b)** FGL; **c)** GGL & Guo et al; and **d)** estimation of time-varying networks (Kolar & Xing, 2009); see Section 2.3 for details.

Unlike proposals of Danaher et al. [6], the estimator of Guo et al. [9] is based on a non-convex penalty, and does not naturally fit into the above framework. However, Lemma 2 in Guo et al. [9] establishes a connection between the optimal solutions of the original optimization problem, with those obtained by considering a single penalty of the form  $\left\{ \sum_{k=1}^K |\theta_{ij}^{(k)}| \right\}^{1/2} \equiv \|\Theta_{ij}\|_{1,2}$ . Similar to GGL, the connection with the method of Guo et al. [9] can be build based on the above alternative formulation, by considering again the empty graph  $G_e$  (Figure 2c), but instead the  $\|\cdot\|_{1,2}$  penalty, which is a member of the CAP family of penalties [36]. More specifically,

$$\left\{ \sum_{k=1}^K |\omega_{ij}^{(k)}| \right\}^{1/2} = \|g_{G_e}(\Theta_{ij})\|_{1,2}.$$

Using the above framework, it is also easy to see the connection between our proposed estimator and the proposal of Kolar et al. [13]: the total variation penalty in Kolar et al. [13] is closely related to FGL, with summation over differences in consecutive time points. It is therefore clear that the penalty of Kolar et al. [13] (up to constant multipliers) can be obtained by applying the graph Laplacian penalty defined for a line graph connecting the time points (Figure 2d).

The above discussion highlights the generality of the proposed estimator, and its connection to existing methods. In particular, while FGL and GGL/Guo et al. [9] consider extreme cases with isolated, or fully connected nodes, one can obtain more flexibility in estimation of multiple precision matrices by defining the penalty based on the known subpopulation network, e.g. based on phylogenetic trees or spatio-temporal similarities between fMRI samples. The clustering-based approach of Section 4 further extends the applicability of the proposed estimator to the settings where the subpopulation network is not known *a priori*. The simulation results in Section 6 show that the additional flexibility of the proposed estimator can result in significant improvements in estimation of multiple precision matrices, when  $K > 2$ . The above discussion also suggests that other variants of the proposed estimator can be defined, by considering other norms. We leave such extensions to future work.

### 3. Theoretical properties

In this section, we establish norm and model selection consistency of the LA-SICH estimator. We consider a high-dimensional setting  $p \gg n_k, k = 1, \dots, K$ , where both  $n$  and  $p$  go to infinity. As mentioned in the Introduction, the normality assumption is not required for establishing these results. We instead require conditions on tails of random vectors  $X^{(k)}$  for each  $k = 1, \dots, K$ . We consider two cases, exponential tails and polynomial tails, which both allow for distributions other than multivariate normal.

**Condition 1** (Exponential Tails). *There exists a constant  $c_1 \in (0, \infty)$  such that*

$$\mathbb{E} \left[ \exp \left\{ t(X_j^{(k)} - \mu_j^{(k)})/(\sigma_{jj}^{(k)})^{1/2} \right\} \right] \leq e^{c_1^2 t^2/2}, \quad \forall t \in \mathbb{R}, k = 1, \dots, K, j = 1, \dots, p.$$

**Condition 2** (Polynomial Tails). *There exist constants  $c_2, c_3 > 0$  and  $c_4$  such that*

$$\mathbb{E} \left[ \left\{ X_j^{(k)}/(\sigma_{jj}^{(k)})^{1/2} \right\}^{4(c_2+c_3+1)} \right] \leq c_4, \quad k = 1, \dots, K, j = 1, \dots, p.$$

Since we adopt the correlation-based Gaussian log-likelihood, we require the boundedness of the true variances to control the error between true and sample correlation matrices.

**Condition 3** (Bounded variance). *There exist constants  $c_5 > 0$  and  $c_6 < \infty$  such that  $c_5 \leq \min_{k,j} \sigma_{jj}^{(k)}$  and  $\max_{k,j} \sigma_{jj}^{(k)} \leq c_6$ .*

**Condition 4** (Sample size). *Let  $\lambda_\Theta \equiv \max_k \|\Theta_0^{(k)}\|_2$ . Let*

$$C_1 \equiv \left\{ 2c_5^{-2} + c_5 + c_6^{-3/2} + 2c_5^{-5/2}c_6 + (c_5^{-4} + 2c_5^{-5}c_6)^{1/2} \right\}^{-1}.$$

(i) (Exponential tails). *It holds that*

$$n \geq \max \left\{ \frac{12}{\min_k \pi_k}, 2^{18} 3^3 C_1^2 (1 + 4c_1^2)^2 c_6^2 \lambda_\Theta^4 \left( 1 + \|L\|_2^{1/2} \right)^2 s \right\} \log p,$$

and  $\log p/n \rightarrow 0$ .

(ii) (Polynomial tails). *Let  $C_2 = \sup_n \{\rho_n \sqrt{n/\log p}\} = O(1)$  where  $\rho_n$  is given in Lemma 1 in the Appendix and  $c_7 > 0$  be some constant. It holds that*

$$n \geq \max \left\{ \frac{p^{1/c_2}}{c_7^{1/c_2}}, 2^7 3^2 C_1^2 C_2^2 K \min_k \pi_k \lambda_\Theta^4 \left( 1 + \|L\|_2^{1/2} \right)^2 s \log p \right\}.$$

Condition 4 determines the sufficient sample size  $n = \sum_k$  for consistent estimation of precision matrices  $\Theta^{(1)}, \dots, \Theta^{(K)}$  in relation to, among other quantities, the number of variables  $p$ , the sparsity pattern  $s$  and the spectral norm of the Laplacian matrix  $\|L\|_2$  of the subpopulation network  $G$ . While a general characterization of  $\|L\|_2$  is difficult, investigating its value in special cases

provides insight into the effect of the underlying population structure on the required sample size. Consider, for instance, two extreme cases: for a fully connected graph  $G$  associated with  $K$  subpopulations,  $\|L\|_2 = 1/(K-1)$ ; for a minimally connected “line” graph, corresponding to e.g. multiple time points,  $\|L\|_2 = 2$ : with  $K = 5$ , 30% more samples are needed for the line graph, compared to a fully connected network. The above calculations match our intuition that fewer samples are needed to consistently estimate precision matrices of  $K$  subpopulations that share greater similarities. This, of course, makes sense, as information can be better shared when estimating parameters of similar subpopulations. Note that, here  $L$  represents the Laplacian matrix of the *true* subpopulation network capturing the underlying population structure. The above conditions thus do not provide any insight into the effect of misspecifying the relationship between subpopulations, i.e., when an incorrect  $L$  is used. This is indeed an important issue that garners additional investigation; see Zhao and Shojaie [37] for some insight in the context of inference for high dimensional regression. In Section 4, we will discuss a data-driven choice of  $L$  that results in consistent estimation of precision matrices.

Before presenting the asymptotic results, we introduce some additional notations. For a matrix  $A = (a_{ij})_{i,j=1}^p \in \mathbb{R}^{p \times p}$ , we denote the spectral norm  $\|A\|_2 = \max_{x \in \mathbb{R}^p, \|x\|=1} \|Ax\|$ , and the element-wise  $\ell_\infty$ -norm  $\|A\|_\infty = \max_{i,j} |a_{i,j}|$  where  $\|x\|$  is the Euclidean norm for a vector  $x$ . We also write the induced  $\ell_\infty$ -norm  $\|A\|_{\infty/\infty} = \sup_{\|x\|_\infty=1} \|Ax\|_\infty$  where  $\|x\|_\infty = \max_i |x_i|$  for  $x = (x_1, \dots, x_p)$ . For the ease of presentation, the results in this section are presented in asymptotic form; non-asymptotic results and proofs are deferred to the Appendix.

### 3.1. Consistency in spectral norm

Let  $s \equiv \#\{(i, j) : \omega_{0,ij}^{(k)} \neq 0, i, j = 1, \dots, p, i \neq j, k = 1, \dots, K\}$ , and  $d = \max_{k,i} \#\{(i, j) : \omega_{0,ij}^{(k)} \neq 0, j = 1, \dots, p, i \neq j\}$ . The following theorem establishes the rate of convergence of the LASICH estimator, in spectral norm, under either exponential or polynomial tail conditions (Condition 1 or 2). Convergence rates for LASICH in  $\ell_\infty$ - and Frobenius norm are discussed in Section 3.3.

**Theorem 1.** *Suppose Conditions 3 and 4 hold. Under Condition 1 or 2,*

$$\sum_{k=1}^K \|\hat{\Omega}_{\rho_n}^{(k)} - \Omega_0^{(k)}\|_2 = O_P \left( \sqrt{\frac{\lambda_{\Theta}^4 (s+1) \log p}{n}} \right),$$

as  $n, p \rightarrow \infty$  where  $\rho_n$  is given in Lemma 1 in the Appendix with  $\gamma = \min_k \pi_k/2$ .

Theorem 1 is proved in the Appendix. The proof builds on tools from Negahban et al. [20]. However, our estimation procedure does not match their general framework: First, we do not penalize the diagonal elements of the inverse correlation matrices; our penalty is thus not a norm. Second, the Laplacian matrix is nonpositive definite. Thus, the Laplacian shrinkage penalty is not strictly convex. The results from Negahban et al. [20] are thus not directly applicable to our

problem. To establish the estimation consistency, we first show, in Lemma 3, that the function  $r(\cdot) = \|\cdot\|_1 + \rho_2 \|\cdot\|_L$  is a seminorm, and is, moreover, convex and decomposable. We also characterize the subdifferential of this seminorm in Lemma 6, based on the spectral decomposition of the graph Laplacian  $L$ . The rest of the proof uses tools from Negahban et al. [20], Rothman et al. [26] and Ravikumar et al. [25], as well as new inequalities and concentration bounds. In particular, in Lemma 4 we establish a new  $\ell_\infty$  bound for the empirical covariance matrix for random variables with polynomial tails, which is used to establish the consistency in the spectral norm under Condition 2.

The convergence rate in Theorem 1 compares favorably to several other methods based on penalized likelihood. Few results are currently available for estimation of multiple precision matrices. An exception is Guo et al. [9], who obtained a slower rate of convergence  $O_p(\{(s+p) \log p/n\}^{1/2})$  under the normality assumption and based on a bound on the Frobenius norm. Our rates of convergence are comparable to the results of Rothman et al. [26] for spectral norm convergence of a single precision matrix, obtained under the normality assumption. Ravikumar et al. [25], on the other hand, assumed the irrepresentability condition to obtain the rate  $O_p(\{\min\{s+p, d^2\} \log p/n\}^{1/2})$  and  $O_p(\{\min\{s+p, d^2\} p^{\tau/(c_2+c_3+1)}/n\}^{1/2})$ , under exponential and polynomial tail conditions, respectively, where  $\tau > 2$  is some scalar. The rate in Theorem 1 is obtained without assuming the irrepresentability condition. In fact, our rates of convergence are faster than those of Ravikumar et al. [25] given the irrepresentability Condition 5 (see Corollary 1). Cai et al. [4] obtained improved rates of convergence under both tail conditions for an estimator that is not found by minimizing the penalized likelihood objective function, and may be nonpositive definite. Finally, note that the results in [26, 25, 4] are for separate estimation of precision matrices and hold for the minimum sample size across subpopulations,  $\min_k n_k$ , whereas our results hold for the total samples size  $\sum_k n_k$ .

### 3.2. Model selection consistency

Let  $S^{(k)} = \{(i, j) : \omega_{0,ij}^{(k)} \neq 0, i, j = 1, \dots, p\}$  be the support of  $\Omega_0^{(k)}$ , and denote by  $d$  the maximum number of nonzero elements in any rows of  $\Omega_0^{(k)}$ ,  $k = 1, \dots, K$ . Define the event

$$\mathcal{M}(\hat{\Omega}_{\rho_n}, \Omega_0) \equiv \left\{ \text{sign}(\hat{\omega}_{\rho_n,ij}^{(k)}) = \text{sign}(\omega_{0,ij}^{(k)}), i, j = 1, \dots, p, k = 1, \dots, K \right\}, \quad (5)$$

where  $\text{sign}(a)$  is 1 if  $a > 0$ , 0 if  $a = 0$  and  $-1$  if  $a < 0$ . We say that an estimator  $\hat{\Omega}_{\rho_n}$  of  $\Omega_0$  is model-selection consistent if  $P\{\mathcal{M}(\hat{\Omega}_{\rho_n}, \Omega_0)\} \rightarrow 1$ .

We begin by discussing an irrepresentability condition for estimation of multiple graphical models. This restrictive condition is commonly assumed to establish model selection consistency of lasso-type estimators, and is known to be almost necessary [19, 35]. For the graphical lasso, Ravikumar et al. [25] showed that the irrepresentability condition amounts to a constraint on the correlation between entries of the Hessian matrix  $\Gamma = \Omega^{-1} \otimes \Omega^{-1}$  in the set  $S$  corresponding to nonzero elements of  $\Omega$ , and those outside this set. Our irrepresentability

condition is motivated by that in Ravikumar et al. [25], however, we adjust the index set  $S$  to also account for covariances of “non-edge variables” that are correlated with each other. More specifically, the description of irrerepresentability condition in Ravikumar et al. [25] involves  $\Gamma_{SS}$  consisting only of elements  $\sigma_{ij}\sigma_{kl}$  with  $(i, j) \in S$  and  $(k, l) \in S$ . However,  $\sigma_{ij} \neq 0$  for  $(i, j) \notin S$  is not taken into account by this definition. We thus adjust the index set  $S$  so that  $\Gamma_{SS}$  also includes elements  $\sigma_{ij}\sigma_{kl}$  if  $(i, k) \in S$  and  $(j, l) \in S$ . This definition is based on the crucial observations that  $\Gamma = \Sigma \otimes \Sigma$  involves the covariance matrix  $\Sigma$  instead of the precision matrix  $\Omega$ , and that some variables are correlated (i.e.,  $\sigma_{ij} \neq 0$ ) even though they may be conditionally independent (i.e.,  $\omega_{ij} = 0$ ). Defining  $S^{(k)}$  for  $k = 1, \dots, K$  as above, we assume the following condition.

**Condition 5** (Irrepresentability condition). *The inverse  $\Theta_0^{(k)}$  of the correlation matrix  $\Psi_0^{(k)}$  satisfies the irrerepresentability condition for  $S^{(k)}$  with parameter  $\alpha$ : (a)  $(\Theta_0^{(k)} \otimes \Theta_0^{(k)})_{S^{(k)}S^{(k)}}$  and  $(\Psi_0^{(k)} \otimes \Psi_0^{(k)})_{S^{(k)}S^{(k)}}$  are invertible, and (b) there exists some  $\alpha \in (0, 1]$  such that*

$$\max_{(i,j) \in (S^{(k)})^c} \|\Gamma_{\{(i,j)\} \times S^{(k)}}^{(k)} \{\Gamma_{S^{(k)}S^{(k)}}^{(k)}\}^{-1}\|_1 \leq 1 - \alpha, \quad (6)$$

for  $k = 1, \dots, K$  where  $\Gamma^{(k)} \equiv \Psi_0^{(k)} \otimes \Psi_0^{(k)}$ .

In addition to the irrerepresentability condition, we require bounds on the magnitude of  $\theta_{ij}^{(k)} \neq 0$  and their normalized difference.

**Condition 6** (Lower bounds for the inverse correlation matrices). *There exists a constant  $c_8 \in \mathbb{R}$  such that*

$$\theta_{\min} \equiv \min_{k=1, \dots, K, i \neq j} |\theta_{0,ij}^{(k)}| \geq c_8 > 0.$$

Moreover, for  $\Omega_{0,ij} \neq 0$ ,  $L\Omega_{0,ij} \neq 0$  and there exists a constant  $c_9 > 0$  such that

$$\min_{l_{kk'} \neq 0, \frac{\omega_{0,ij}^{(k)}}{\sqrt{d_k}} - \frac{\omega_{0,ij}^{(k')}}{\sqrt{d_{k'}}} \neq 0} \left| \frac{\theta_{0,ij}^{(k)}}{\sqrt{d_k}} - \frac{\theta_{0,ij}^{(k')}}{\sqrt{d_{k'}}} \right| \geq c_9.$$

The first lower bound in Condition 6 is the usual “min-beta” condition for model selection consistency of lasso-type estimators. The second lower bound, which is represented here for the normalized Laplacian penalty, is a mild condition which ensures estimates based on inverse correlation matrices can be mapped to precision matrices. For any pair of subpopulations  $k$  and  $k'$  connected in  $G$  it requires that if the difference in (normalized) entries of the entires of the precision matrices are nonzero, the difference in (normalized) entries of inverse correlation matrices are bounded away from zero. In other words, the bound guarantees that  $\Theta_{0,ij}$  is not in the null space of  $L$ , whenever  $\Omega_{0,ij}$  is outside of the null space. This bound can be relaxed if we use a positive definite matrix  $L_\epsilon = L + \epsilon I$  for  $\epsilon > 0$  small.

Our last condition for establishing the model selection consistency concerns the minimum sample size and the tuning parameter for the graph Laplacian penalty. This condition is necessary to control the  $\ell_\infty$ -bound of the error  $\hat{\Theta}_{\rho_n} - \Theta_0$ , as in Ravikumar et al. [25]. Our minimum sample size requirement is related to the irrepresentability condition. Let  $\kappa_\Gamma$  be the maximum of the absolute column sums of the matrices  $\{(\Gamma^{(k)})^{-1}\}_{S^{(k)}S^{(k)}}, k = 1, \dots, K$ , and  $\kappa_\Psi$  be the maximum of the absolute column sums of the matrices  $\Psi_0^{(k)}, k = 1, \dots, K$ . The minimum sample size in Ravikumar et al. [25] is also a function of the irrepresentability constant, in particular, their  $\kappa_\Gamma$  involves  $\{(\Gamma_{S^{(k)}S^{(k)}}^{(k)})\}^{-1}$ . There is, therefore, a subtle difference between our definition and theirs: in our definition, the matrix is first inverted and then partitioned, while in Ravikumar et al. [25], the matrix is first partitioned and then inverted. Corollary 2 establishes the model selection consistency under a weaker sample size requirement, by exploiting instead the control of the spectral norm in Theorem 1.

**Condition 7** (Sample size and regularization parameters). *Let*

$$C_3 = \max \left\{ \frac{2^6 3^4 \kappa_\Psi^2 \kappa_\Gamma^2}{\min_k \pi_k^2} \max \left\{ 1, \frac{2^6 7^2 \kappa_\Psi^4 \kappa_\Gamma^2}{\alpha^2 \min_k \pi_k^2} \right\}, \frac{36}{c_8^2}, \frac{2^4 3^2}{c_9^2 \min_k d_k} \right\}$$

(i) (Exponential tails). *It holds*

$$n > \frac{12 \log p}{\min_k \pi_k} \max \{1, 2^6 3^2 C_1^2 (1 + c_1^2)^2 c_6^2 C_3 d^2\}$$

(ii) (Polynomial tails). *It holds  $n > \max\{p^{1/c_2} c_7^{-1/c_2}, C_1^2 C_2^2 C_3 d^2 \log p\}$ .*

(iii) *It holds that  $\rho_2 \leq \alpha^2 / \{4 \|L\|_2^{1/2} (2 - \alpha)\}$ .*

With these condition, we obtain

**Theorem 2.** *Suppose that Conditions 3, 5, 6 and 7 hold. Under Condition 1 or 2,  $P(\mathcal{M}(\hat{\Omega}_{\rho_n}, \Omega_0)) \rightarrow 1$  as  $n, p \rightarrow \infty$  where  $\rho_n$  is given in Lemma 1 in the Appendix with  $\gamma = \min_k \pi_k / 2$ .*

### 3.3. Additional results

In this section, we establish norm and variable selection consistency of LASICH under alternative assumptions. Our first result gives better rates of convergence for consistency in the  $\ell_\infty$ -, spectral and Frobenius norms, under the condition for model selection consistency. Our rates in Corollary 1 improve the previous results by Ravikumar et al. [25], and are comparable to that of Cai et al. [4] in the  $\ell_\infty$ - and spectral norms under both tail conditions.

**Corollary 1.** *Suppose the conditions in Theorem 2 hold. Then, under Condition 1 or 2,*

$$\sum_{k=1}^K \|\hat{\Omega}_{\rho_n}^{(k)} - \Omega_0^{(k)}\|_F = O_P \left( \sqrt{\frac{\min\{\lambda_{\Theta}^4 p(s+1), \kappa_\Gamma^2 (s+p)\} \log p}{n}} \right),$$

$$\sum_{k=1}^K \|\hat{\Omega}_{\rho_n}^{(k)} - \Omega_0^{(k)}\|_2 = O_P \left( \sqrt{\frac{\min\{\lambda_{\Theta}^4(s+1), \kappa_{\Gamma}^2 d^2\} \log p}{n}} \right),$$

$$\sum_{k=1}^K \|\hat{\Omega}_{\rho_n}^{(k)} - \Omega_0^{(k)}\|_{\infty} = O_P \left( \sqrt{\frac{\kappa_{\Gamma}^2 \log p}{n}} \right).$$

Our next result in Corollary 2 establishes the model selection consistency under a weaker version of the irrepresentability condition (Condition 6). Aside from the difference in the index sets  $S^{(k)}$ , the form of the Condition 6 and the assumption of invertibility of  $(\Psi_0^{(k)} \otimes \Psi_0^{(k)})_{S^{(k)} S^{(k)}}$  are similar to those in Ravikumar et al. [25]. On the other hand, Ravikumar et al. [25] do not require invertibility of  $(\Theta_0^{(k)} \otimes \Theta_0^{(k)})_{S^{(k)} S^{(k)}}$ . However, their proof is based on an application of Brouwer's fixed point theorem, which does not hold for the corresponding function (Eq. (70) in page 973) since it involves a matrix inverse, and is hence not continuous on its range. The additional inevitability assumption in Condition 6 is used to address this issue in Lemma 11. The condition can be relaxed if we assume an alternative scaling of the sample size stated in Condition 8 below instead of Condition 7.

**Condition 8.** Let  $\lambda_{\Psi} = \max_k \|\Psi_0^{(k)}\|$ . Suppose  $\rho_2 \leq \alpha^2 / \{4\|L\|_2^{1/2}(2 - \alpha)\}$  and (i) (Exponential tails)

$$n > 2^{19} 3^3 \{\min_k \pi_k\}^{-3} C_1^2 (1 + 4c_1^2)^2 c_6^2 \lambda_{\Theta}^4 \left(1 + \rho_2 \|L\|_2^{1/2}\right)^2 s \log p \max\{\lambda_{\Psi}, 4\lambda_{\Theta}^4 \alpha^{-1}\},$$

or

(ii) (Polynomial tails)

$$n > 2^{12} 3^3 \{\min_k \pi_k\}^{-2} K^2 C_1^2 C_2^2 \lambda_{\Theta}^4 \left(1 + \rho_2 \|L\|_2^{1/2}\right)^2 s \log p \max\{\lambda_{\Psi}, 4\lambda_{\Theta}^4 \alpha^{-1}\}.$$

**Corollary 2.** Suppose that Conditions 3, 6 and 8 hold. Suppose also that Condition 5 holds without requiring the invertibility of  $(\Theta_0^{(k)} \otimes \Theta_0^{(k)})_{S^{(k)} S^{(k)}}$ . Then, under Condition 1 or 2,  $P(\mathcal{M}(\hat{\Omega}_{\rho_n}, \Omega_0)) \rightarrow 1$  as  $n, p \rightarrow \infty$  where  $\rho_n$  is given in Lemma 1 in the Appendix with  $\gamma = \min_k \pi_k / 2$ .

#### 4. Laplacian shrinkage based on hierarchical clustering

Our proposed LASICH approach utilizes the information in the subpopulation network  $G$ . In practice, however, similarity between subpopulations may be difficult to ascertain or quantify. In this section, we present a modified LASICH framework, called HC-LASICH, which utilizes hierarchical clustering to learn the relationships among subpopulations. The information from hierarchical clustering is then used to define the weighted subpopulation network. Importantly, HC-LASICH can even be used in settings where the subpopulation membership is unavailable, for instance, to learn the genetic network of cancer patients, where cancer subtypes may be unknown.

We use hierarchical clustering with a complete, single or average linkage to estimate both the subpopulation memberships and the weighted subpopulation network  $G$ . Specifically, the length of a path between two subpopulations in the dendrogram is used as a measure of dissimilarity between two subpopulations; the weights for the subpopulation networks are simply defined by taking the inverse of these lengths. Throughout this section, we assume that the number of subpopulations  $K$  is known. While a number of methods have been proposed for estimating the number of subpopulations in hierarchical clustering (see e.g. Borysov et al. [1] and the references therein), the problem is beyond the scope of this paper.

Let  $\mathcal{I} = (\mathcal{I}^{(1)}, \dots, \mathcal{I}^{(K)})$  be the subpopulation membership indicator such that  $\mathcal{I}$  follows the multinomial distribution  $\text{Mult}_K(1, (\pi_1, \dots, \pi_K))$  with parameter 1 and subpopulation membership probabilities  $(\pi_1, \dots, \pi_K) \in (0, 1)^K$ . Note that  $\mathcal{I}$  is missing and is to be estimated. Let  $\mathcal{I}_i, i = 1, \dots, n$  be i.i.d. copies of  $\mathcal{I}$  and  $\hat{\mathcal{I}}_i = (\hat{\mathcal{I}}_i^1, \dots, \hat{\mathcal{I}}_i^K)$  be an estimated subpopulation indicator for the  $i$ th observation via hierarchical clustering. Based on the estimated subpopulation membership and subpopulation network  $\hat{G}$ , we apply our method to obtain the estimator, HC-LASICH,  $\hat{\Omega}_{HC, \rho_n} = (\hat{\Omega}_{HC, \rho_n}^{(1)}, \dots, \hat{\Omega}_{HC, \rho_n}^{(K)})$ . Interestingly, HC-LASICH enjoys the same theoretical properties as LASICH, under the normality assumption. To show this, we first establish the consistency of hierarchical clustering in high dimensions, which is of independent interest. Our result is motivated by the recent work of [1], who study the consistency of hierarchical clustering for independent normal variables  $X^{(k)} \sim N(\mu^{(k)}, \sigma^{(k)}I)$ ; we establish similar results for multivariate normal distributions with arbitrary covariance structures. We make the following assumption.

**Condition 9.** For  $k, k' = 1, \dots, K$ , let

$$\bar{\lambda}^{(k)} = p^{-1} \sum_{j=1}^p \lambda^{(k),j},$$

$$\mu^{(k,k')} = p^{-1} \left\| \Lambda_{k,k'}^{1/2} Q_{k,k'}^T \left[ \Sigma^{(k)} + \Sigma^{(k')} \right]^{1/2} \left[ \mu^{(k)} - \mu^{(k')} \right] \right\|^2,$$

where  $\lambda^{(k),j}$  is the eigenvalues of  $\Sigma^{(k)}$  with  $\lambda^{(k),1} \leq \lambda^{(k),2} \leq \dots \leq \lambda^{(k),p}$ , and the spectral decomposition of  $\Sigma^{(k)} + \Sigma^{(k')}$  is  $\Sigma^{(k)} + \Sigma^{(k')} = Q_{k,k'} \Lambda_{k,k'} Q_{k,k'}^T$ . It holds that

$$\mu^{(k,k')} > 2 \min \left\{ \bar{\lambda}^{(k)}, \bar{\lambda}^{(k')} \right\} - \lambda^{(k),p} - \lambda^{(k'),p}, \quad k \neq k', \quad k, k' = 1, \dots, K,$$

$$0 < c_{10} \leq \lambda^{(k),j} \leq c_{11} < \infty, \quad \|\mu^{(k)}\| \leq c_{11}, \quad k = 1, \dots, K, j = 1, \dots, p.$$

for constants  $m$  and  $M$ .

Under the normality assumption, the following results shows that the probability of successful clustering converges to 1, as  $p, n \rightarrow \infty$ .



**Theorem 3.** Suppose that that  $X^{(k)}, k = 1, \dots, K$ , is normally distributed. Under Condition 9,

$$P(\hat{\mathcal{I}}_i = \mathcal{I}_i, i = 1, \dots, n) \rightarrow 1, \quad (7)$$

as  $n, p \rightarrow \infty$ .

To proof of Theorem 3 generalizes recent results of Borysov et al. [1] to the case of arbitrary covariance structures. A key component of the proof is a new bound on the  $\ell_2$  norm of a multivariate normal random variable with arbitrary mean and covariance matrix established in Lemma 14. The proof of the lemma uses new concentration inequalities for high-dimensional problems in [2], and may be of independent interest.

Note that the consistent estimation of subpopulation memberships (7) implies that the estimated hierarchy among clusters also matches the true hierarchy. Thus, with successful clustering established in Theorem 3, theoretical properties of  $\hat{\Omega}_{HC, \rho_n}$  naturally follow.

**Theorem 4.** Suppose that  $X^{(k)}, k = 1, \dots, K$ , is normally distributed and that Condition 9 holds. (i) Under the conditions of Theorem 1,

$$\sum_{k=1}^K \|\hat{\Omega}_{HC, \rho_n}^{(k)} - \Omega_0^{(k)}\|_2 = O_P \left( \sqrt{\frac{\lambda_{\Theta}^4 (s+1) \log p}{n}} \right).$$

Suppose, moreover, that the conditions of Theorem 2 holds. Then

$$\begin{aligned} \sum_{k=1}^K \|\hat{\Omega}_{HC, \rho_n}^{(k)} - \Omega_0^{(k)}\|_F &= O_P \left( \sqrt{\frac{\min\{\lambda_{\Theta}^4 p (s+1), \kappa_{\Gamma}^2 (s+p)\} \log p}{n}} \right), \\ \sum_{k=1}^K \|\hat{\Omega}_{HC, \rho_n}^{(k)} - \Omega_0^{(k)}\|_2 &= O_P \left( \sqrt{\frac{\min\{\lambda_{\Theta}^4 (s+1), \kappa_{\Gamma}^2 d^2\} \log p}{n}} \right), \\ \sum_{k=1}^K \|\hat{\Omega}_{HC, \rho_n}^{(k)} - \Omega_0^{(k)}\|_{\infty} &= O_P \left( \sqrt{\frac{\kappa_{\Gamma}^2 \log p}{n}} \right). \end{aligned}$$

(ii) Under the conditions of Theorem 2,

$$P(\mathcal{M}(\hat{\Omega}_{HC, \rho_n}, \Omega_0)) \rightarrow 1, \quad \text{as } n, p \rightarrow \infty.$$

## 5. Algorithms

We develop an alternating directions method of multipliers (ADMM) to efficiently solve the convex optimization problem (3).

Let  $A^{(k)} = (a_{ij}^{(k)})_{i,j=1}^p \in \mathbb{R}^{p \times p}$ ,  $B^{(k)} = (b_{ij}^{(k)})_{i,j=1}^p \in \mathbb{R}^{p \times p}$ ,  $C^{(k)} = (c_{ij}^{(k)})_{i,j=1}^p \in \mathbb{R}^{p \times p}$ ,  $D^{(k)} = (d_{ij}^{(k)})_{i,j=1}^p \in \mathbb{R}^{p \times p}$ ,  $k = 1, \dots, K$ . Define  $A = (A^{(1)}, \dots, A^{(K)})$ ,  $B = (B^{(1)}, \dots, B^{(K)})$ ,  $C = (C^{(1)}, \dots, C^{(K)})$ ,  $D = (D^{(1)}, \dots, D^{(K)})$ , and  $c_{ij} \equiv (c_{ij}^{(1)}, \dots, c_{ij}^{(K)})^T \in \mathbb{R}^K$ ,  $d_{ij} \equiv (d_{ij}^{(1)}, \dots, d_{ij}^{(K)})^T \in \mathbb{R}^K$ ,  $e_{C,ij} \equiv (e_{C,ij}^{(1)}, \dots, e_{C,ij}^{(K)})^T \in \mathbb{R}^K$  where  $E_C^{(k)} = (e_{C,ij}^{(k)})_{i,j=1}^p$ .

To facilitate the computation, we consider instead a perturbed graph Laplacian  $L_\epsilon = L + \epsilon I$ , where  $I$  is the identity matrix and  $\epsilon > 0$  is a small perturbation. The difference between solutions to the original and modified optimization problem is largely negligible for small  $\epsilon$ ; however, the positive definiteness of  $L_\epsilon$  results in more efficient computation. A similar idea was used in Guo et al. [9] and Rothman et al. [26] to avoid division by zero. The optimization problem (3) with  $L$  replaced by  $L_\epsilon$  can then be written as

$$\begin{aligned} & \text{minimize} \quad \sum_{k=1}^K \frac{n_k}{n} \left( \text{tr} \left( \Psi_n^{(k)} A^{(k)} \right) - \log \det(A^{(k)}) \right) + \rho_n \sum_{k=1}^K \|B^{(k)}\|_1 \\ & \quad + \rho_n \rho_2 \sum_{i \neq j} (c_{ij}^T L_\epsilon c_{ij})^{1/2} \\ & \text{s.t.} \quad A^{(k)} = D^{(k)}, B^{(k)} = D^{(k)}, L_\epsilon c_{ij} = L_\epsilon d_{ij} \quad k = 1, \dots, K, i, j = 1, \dots, p. \end{aligned} \quad (8)$$

Using Lagrange multipliers  $E = (E_A, E_B, E_C)^T$ , with  $E_A = (E_A^{(1)}, \dots, E_A^{(K)})$  with  $E_A^{(k)} \in \mathbb{R}^{p \times p}$ ,  $k = 1, \dots, K$ ,  $E_B = (E_B^{(1)}, \dots, E_B^{(K)})$  with  $E_B^{(k)} \in \mathbb{R}^{p \times p}$ ,  $k = 1, \dots, K$ , and  $E_C = (E_C^{(1)}, \dots, E_C^{(K)})$  with  $E_C^{(k)} \in \mathbb{R}^{p \times p}$ ,  $k = 1, \dots, K$ , the augmented Lagrangian in scaled form is given by

$$\begin{aligned} & L_\varrho(A, B, C, D, E) \\ & \equiv n^{-1} \sum_{k=1}^K n_k \left( \text{tr} \left( \Psi_n^{(k)} A^{(k)} \right) - \log \det(A^{(k)}) \right) + \rho_n \sum_{k=1}^K \|B^{(k)}\|_1 \\ & \quad + \rho_n \rho_2 \sum_{i \neq j} (c_{ij}^T L_\epsilon c_{ij})^{1/2} + \frac{\varrho}{2} \sum_{k=1}^K \left\| A^{(k)} - D^{(k)} + E_A^{(k)} \right\|_F^2 \\ & \quad + \frac{\varrho}{2} \sum_{k=1}^K \left\| B^{(k)} - D^{(k)} + E_B^{(k)} \right\|_F^2 + \frac{\varrho}{2} \sum_{i,j} \left\| L_\epsilon^{1/2} c_{ij} - L_\epsilon^{1/2} d_{ij} + e_{C,ij} \right\|_F^2. \end{aligned}$$

Here  $\varrho > 0$  is a regularization parameter and  $L_\epsilon^{1/2}$  is the square root of  $L_\epsilon$  with  $L_\epsilon = (L_\epsilon^{1/2})^T L_\epsilon^{1/2}$ .

The proposed ADMM algorithm is as follows.

- *Step 0.* Initialize  $A^{(k)} = A^{(k),0}$ ,  $B^{(k)} = B^{(k),0}$ ,  $C^{(k)} = C^{(k),0}$ ,  $D^{(k)} = D^{(k),0}$ ,  $E_A^{(k)} = E_A^{(k),0}$ ,  $E_B^{(k)} = E_B^{(k),0}$ ,  $E_C^{(k)} = E_C^{(k),0}$  and choose  $\varrho > 0$ . Select a scalar  $\varrho > 0$ .
- *Step m.* Given the  $(m-1)$ th estimates,
  - Update  $A^{(k)}$  Find  $A^m$  minimizing  $-\ell_n(A) - (\varrho/2) \sum_{k=1}^K \|A^{(k)} - D^{(k),m-1} - E_A^{(k),m-1}\|$  (see pages 46-47 of Boyd et al. [3] for details).
  - (Update  $B^{(k)}$ ) Compute  $B_{ij}^{(k),m} = S_{\rho_n/\varrho}(D_{ij}^{(k),m-1} - E_{B,ij}^{(k),m-1})$ , where  $S_y(x)$  is  $x - y$  if  $x > y$ , is 0 if  $|x| \leq y$ , and is  $x + y$  if  $x < -y$ .

– (Update  $C^{(k)}$ ) For  $(x)_+ = \max\{x, 0\}$ , compute

$$c_{ij}^m = \left( 1 - \frac{\rho_n \rho_2}{\varrho \|L_\epsilon^{1/2} d_{ij}^{m-1} - e_{C,ij}^{m-1}\|} \right)_+ (d_{ij}^{m-1} - L_\epsilon^{-1/2} e_{C,ij}^{m-1}).$$

– (Update  $D^{(k)}$ ) Compute

$$d_{ij}^m = (2I + L_\epsilon)^{-1} \{a_{ij}^m + e_{A,ij}^{m-1} + b_{ij}^m + e_{B,ij}^{m-1} + L_\epsilon c_{ij}^m + (L_\epsilon^{1/2})^T e_{C,ij}^{m-1}\}.$$

– (Update  $E_A$ ) Compute  $E_A^{(k),m} = E_A^{(k)} + A^{(k),m} - D^{(k),m}$ .

– (Update  $E_B$ ) Compute  $E_B^{(k),m} = E_B^{(k)} + B^{(k),m} - D^{(k),m}$ ,

– (Update  $E_C$ ) Compute  $e_{C,ij}^{(k),m} = e_{C,ij}^{(k)} + L^{1/2}(c_{ij}^{(k),m} - d_{ij}^{(k),m})$ .

- Repeat the iteration until the maximum of the errors  $r_A^{(k)} = A^{(k)} - D^{(k)}$ ,  $r_B^{(k),m} = B^{(k),m} - D^{(k),m}$ ,  $r_C^{(k),m} = C^{(k),m} - D^{(k),m}$ ,  $s^{(k),m} = \varrho(D^{(k),m} - D^{(k),m-1})$  in the Frobenius norm is less than a specified tolerance level.

The proposed ADMM algorithm facilitates the estimation of parameters of moderately large problems. However, parameter estimation in high dimensions can be computationally challenging. We next present a result that determines whether the solution to the optimization problem (3), for given values of tuning parameters  $\rho_n, \rho_2$ , is block diagonal. (Note that this result is an *exact* statement about the *solution* to (3), and does not assume block sparsity of the true precision matrices; see Theorems 1 and 2 of Danaher et al. [6] for similar results.) More specifically, the condition in Proposition 1 provides a very fast check, based on the entries of the empirical correlation matrices  $\Psi_n^{(k)}, k = 1, \dots, K$ , to identify the block sparsity pattern in  $\hat{\Omega}_{\rho_n}^{(k)}, k = 1, \dots, K$  after some permutation of the features.

Let  $U_L = [u_1 \dots u_K] \in \mathbb{R}^{K \times K}$  where  $u_1, \dots, u_K$ 's are eigenvectors of  $L$  corresponding to  $0, \lambda_{L,2}, \dots, \lambda_{L,K}$ . Define  $\Lambda_L^{-1/2}$  as the diagonal matrix with diagonal elements  $0, \lambda_{L,2}^{-1/2}, \dots, \lambda_{L,K}^{-1/2}$ .

**Proposition 1.** *The solution  $\hat{\Omega}_{\rho_n}^{(k)}, k = 1, \dots, K$  to the optimization problem (3) consists of the block diagonal matrices with the same block structure  $\text{diag}(\Omega_1, \dots, \Omega_B)$  among all groups if and only if for  $\Psi_{n,ij} = (\psi_{n,ij}^{(1)}, \dots, \psi_{n,ij}^{(K)})^T$*

$$\min_{v \in [-1,1]^K} \left\| \Lambda_L^{-1/2} U_L \left( \frac{n_k}{n} \Psi_{n,ij} - \rho_n v \right) \right\| \leq \rho_n \rho_2, \quad (9)$$

and for all  $i, j$  such that the  $(i, j)$  element is outside the blocks.

The proof of the Proposition is similar to Theorems 1 of Danaher et al. [6] and is hence omitted. Condition 9 can be easily verified by applying quadratic programming to the left hand side of the inequality. The solution to (3) can then be equivalently found by solving the optimization problem separately for each of the blocks; this can result in significant computational advantages for moderate to large values of  $\rho_n \rho_2$ .

## 6. Numerical results

### 6.1. Simulation experiments

We compare our method with four existing methods, graphical lasso, the method of Guo et al. [9], FGL and GGL of Danaher et al. [6]. For graphical lasso, estimation was carried out separately for each group with the same regularization parameter.

Our simulation setting is motivated by estimation of gene networks for healthy subjects and patients with two similar diseases caused by inactivation of certain biological pathways. We consider  $K = 3$  groups with sample sizes  $n = (50, 100, 50)$  and dimension  $p = 100$ . Data are generated from multivariate normal distributions  $N(\mu^{(k)}, (\Omega_0^{(k)})^{-1})$ ,  $k = 1, 2, 3$ ; all precision matrices  $\Omega_0^{(k)}$  are block diagonal with 4 blocks of equal size.

To create the precision matrices, we first generated a graph with 4 components of equal size, each as either an Erdős-Rényi or scale free graphs with 95 total edges. We randomly assigned  $\text{Unif}((-0.7, -0.5) \cup (0.5, 0.7))$  values to nonzero entries of the corresponding adjacency matrix  $A$  and obtained a matrix  $\tilde{A}$ . We then added 0.1 to the diagonal of  $\tilde{A}$  to obtain a positive definite matrix  $\Omega_0^{(1)}$ . For each of subpopulations 2 and 3, we removed one of the components of the graph by setting the off diagonal entries of  $\tilde{A}$  to zero, and added a perturbation from  $\text{Unif}(-0.2, 0.2)$  to nonzero entries in  $\tilde{A}$ . Positive definite matrices  $\Omega_0^{(2)}$  and  $\Omega_0^{(3)}$  were obtained by adding 0.1 to the diagonal elements. All partial correlations ranges from .28 to .54 in the absolute values. A similar setting was considered in Danaher et al. [6], where the graph included more components, but no perturbation was added. We consider two simulation settings, with *known* and *unknown* subpopulation network  $G$ .

#### 6.1.1. Known subpopulation network $G$

In this case, we set  $\mu^{(k)} = 0$ ,  $k = 1, 2, 3$  and use the graph in Figure 1 as the subpopulation network.

Figures 3a,c show the average number of true positive edges versus the average number of detected edges over 50 simulated data sets. Results for multiple choices of the second tuning parameter are presented for FGL, GGL and LASICH. It can be seen that in both cases, LASICH outperforms other methods, when using relatively large values of  $\rho_2$ . Smaller values of  $\rho_2$ , on the other hand, give similar results as other methods of joint estimation of multiple graphical models. These results indicate that, when the available subpopulation network is informative, the Laplacian shrinkage constraint can result in significant improvement in estimation of the underlying network.

Figures 3b,d show the estimation error, in Frobenius norm, versus the number of detected edges. LASICH has larger errors when the estimated graphs have very few edges, but, its error decreases as the number of detected edges increase, eventually yielding smaller errors than other methods. The non-convex penalty

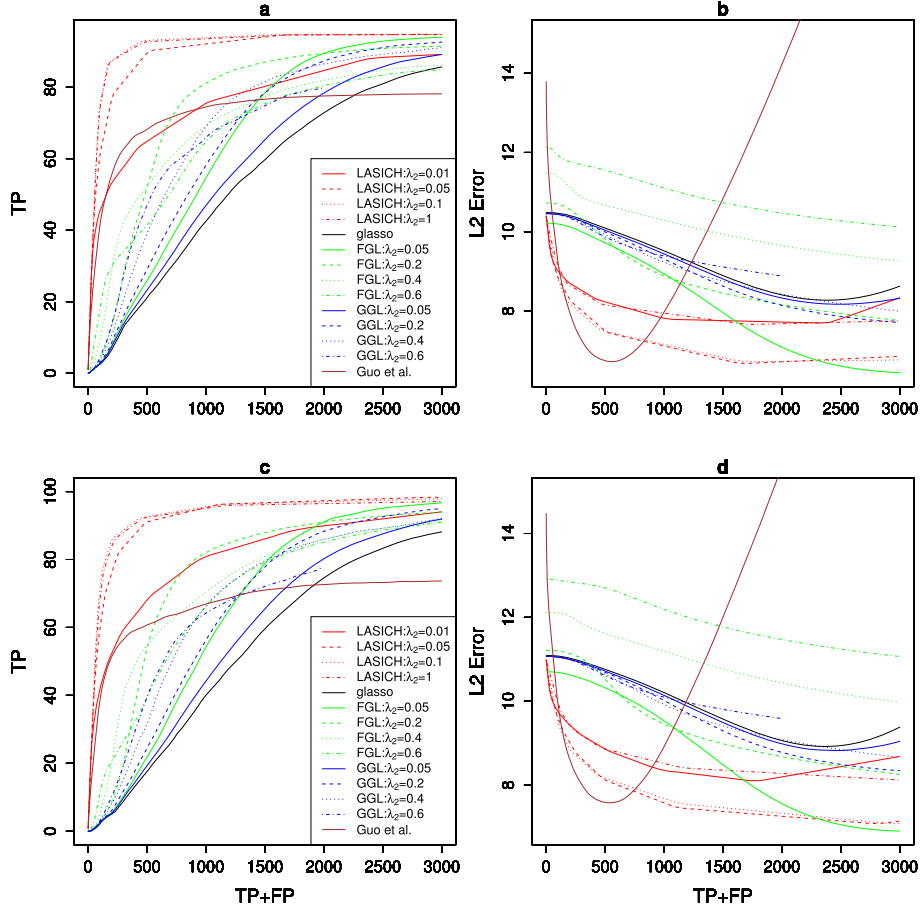


FIG 3. Simulation results for joint estimation of multiple precision matrices with known subpopulation memberships. Results show the average number of true positive edges (a & c) and estimation error, in Frobenius norm (b & d) over 50 data sets with  $n = 200$  multivariate normal observations generated from a graphical model with  $p = 100$  features; results in top row (a & b) are for an Erdős-Rényi graph and those in bottom row (c & d) are for a scale free (power-law) graph.

of Guo et al. [9] performs well in terms of estimation error, although determining the appropriate range of tuning parameter for this method may be difficult.

#### 6.1.2. Unknown subpopulation network $G$

In this case, the subpopulation memberships and the subpopulation network  $G$  are estimated based on hierarchical clustering. We randomly generated  $\mu^{(1)}$  from a multivariate normal distribution with a covariance matrix  $\sigma^2 I$ . For subpopulations 2 and 3, the elements of  $\mu^{(1)}$  corresponding to the empty components of

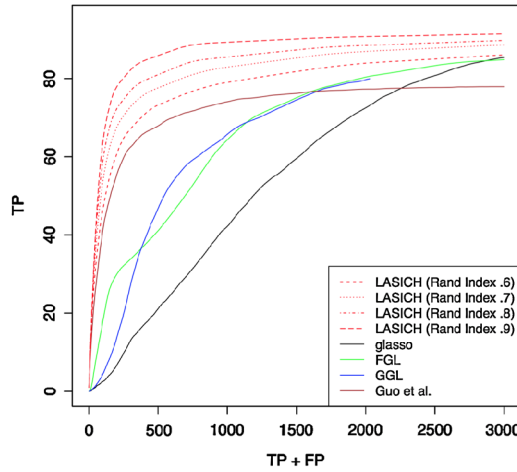


FIG 4. Simulation results for joint estimation of multiple precision matrices with unknown subpopulation memberships. Results show the average number of true positive edges over 50 data sets with  $n = 200$  multivariate normal observations generated from a graphical model with over an Erdős-Rényi graph with  $p = 100$  features. Results for HC-LASICH and FGL/GGL correspond to the best choice of the second tuning parameter among those in Figure 3a. The Rand indices for HC-LASICH are averages over 50 generated data sets.

the graph were set to zero to obtain  $\mu^{(2)}$  and  $\mu^{(3)}$ . Hierarchical clustering with complete linkage was applied to data to obtain the dendrogram; we took inverse of distances in the dendrogram to obtain similarity weights used in the graph Laplacian.

Figure 4 compares the performance of HC-LASICH, in terms of support recovery, to competing methods, in the setting where the subpopulation memberships and network are estimated from data (Section 4). Here the differences in subpopulation means  $\mu^{(k,k')}$  are set up to evaluate the effect of clustering accuracy. The four settings considered correspond to average Rand indices of .6 .7, .8 and .9 across 50 data sets, respectively. Here the second tuning parameter for HC-LASICH, GGL and FGL is chosen according to the best performing model in Figure 3. As expected, changing the mean structure, and correspondingly the Rand index, does not affect the performance of other methods. The results indicate that, as long as features can be clustered in a meaningful way, HC-LASICH can result in improved support recovery. Data-adaptive choices of the tuning parameter corresponding to the Laplacian shrinkage penalty may result in further improvements in the performance of the HC-LASICH. However, we do not pursue such choices here.

## 6.2. Genetic networks of cancer subtypes

Breast cancer is heterogenous with multiple clinically verified subtypes [22]. Jönsson et al. [12] used copy number variation and gene expression measure-

ments to identify new subtypes of breast cancer and showed that the identified subtypes have distinct clinical outcomes. The genetic networks of these different subtypes are expected to share similarities, but to also have unique features. Moreover, the similarities among the networks are expected to corroborate with the clustering of the subtypes based on their molecular profiles. We applied network estimation methods of Section 6.1 to a subset of the microarray gene expression data from Jönsson et al. [12], containing data for 218 patients classified into three previously known subtypes of breast cancer: 46 Luminal-simple, 105 Luminal-complex and 67 Basal-complex samples. For ease of presentation, we focused on 50 genes with largest variances. The hierarchical clustering results of Jönsson et al. [12], reproduced in Figure 5 for the above three subtypes, were used to identify the subpopulation membership; reciprocals of distances in the dendrogram were used to define similarities among subtypes used in the graph Laplacian penalty.

To facilitate the comparison, tuning parameters were selected such that the estimated networks of the three subtypes using each method contained a total of 150 edges. For methods with two tuning parameters, pairs of tuning parameters were determined using the Bayesian information criterion (BIC), as described in Guo et al. [9]. Estimated genetic networks of the three cancer subtypes are shown in Figure 5. For each method, edges common in all three subtypes, those common in Luminal subtypes and subtype specific edges are distinguished.

In this example, separate graphical lasso estimates and FGL/GGL estimates are two extremes. Estimated network topologies from graphical lasso vary from subtype to subtype, and common structures are obscured; this variability may be because similarities among subtypes are not incorporated in the estimation. In contrast, FGL and GGL give identical networks for all subtypes, perhaps because both methods encourage the estimated networks of all subtypes to be equally similar. Intermediate results are obtained using LASICH and the method of Guo et al. [9]. The main difference between these two methods is that Guo et al. [9] finds more edges common to all three subtypes, whereas LASICH finds more edges common to the Luminal subtypes. This difference is likely because LASICH prioritizes the similarity between the Luminal subtypes via graph Laplacian while the method of Guo et al. [9] does not distinguish between the three subtypes. The above example highlights the potential advantages of LASICH in providing network estimates that better corroborate with the known hierarchy of subpopulations.

## 7. Discussion

We introduced a flexible method for joint estimation of multiple precision matrices, called LASICH, which is particularly suited for settings where observations belong to three or more subpopulations. In the proposed method, the relationships among heterogeneous subpopulations is captured by a weighted network, whose nodes correspond to subpopulations, and whose edges capture their similarities. As a result, LASICH can model complex relationships among subpopulations, defined, for example, based on hierarchical clustering of samples.

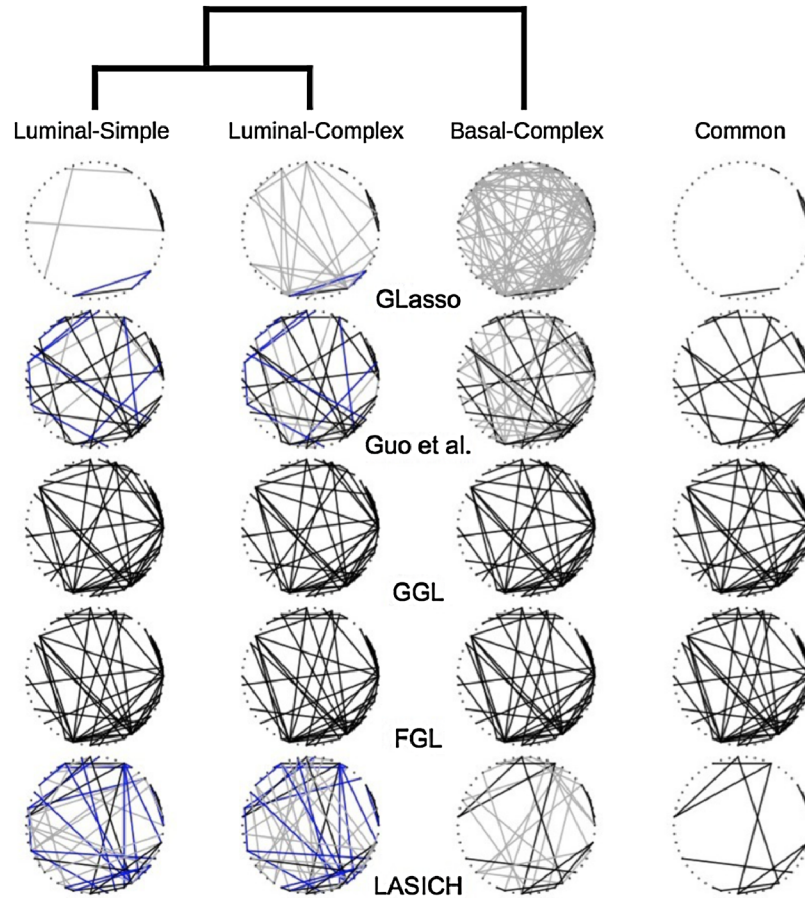


FIG 5. Dendrogram of hierarchical clustering of three subtypes of breast cancer from Jönsson et al. (2010) along with estimated gene networks using graphical lasso (GLasso), method of Guo et al., FGL and GGL of Daneher et al. (2014) and LASICH. Blue edges are common to Luminal subtypes and black edges are shared by all three subtypes; condition specific edges are drawn in gray.

We established asymptotic properties of the proposed estimator in the setting where the relationship among subpopulations is externally defined. We also extended the method to the setting of unknown relationships among subpopulations, by showing that clusters estimated from the data can accurately capture the true relationships. The proposed method generalizes existing convex penalties for joint estimation of graphical models, and can be particularly advantageous in settings with multiple subpopulations.

A particularly appealing feature of the proposed extension of LASICH is that it can also be applied in settings where the subpopulation memberships are unknown. The latter setting is closely related to estimation of precision matrices for mixture of Gaussian distributions. Both approaches have limitations and draw-



backs: on the one hand, the extension of LASICH to unknown subpopulation memberships requires certain assumptions on differences of population means (Section 4). On the other hand, estimation of precision matrices for mixture of Gaussians is computationally challenging, and known rates of convergence of parameter estimation in mixture distributions (e.g. in Städler et al. [29]) are considerably slower.

Throughout this paper we assumed that the number of subpopulations is known. Extensions of this method to estimation of graphical models in populations with an unknown number of subpopulations would be particularly interesting for analysis of genetic networks associated with heterogeneity in cancer samples, and are left for future research.

## 8. Appendix: Proofs and technical details

We denote true inverse correlation matrices as  $\Theta_0 = (\Theta_0^{(1)}, \dots, \Theta_0^{(K)})$  and true correlation matrices as  $\Psi_0 = (\Psi_0^{(1)}, \dots, \Psi_0^{(K)})$ , where  $\Theta_0^{(k)} \equiv (\Psi_0^{(k)})^{-1} \equiv (\theta_{0,ij}^{(k)})_{i,j=1}^p$ , and  $\Psi_0^{(k)} = (\psi_{0,ij}^{(k)})_{i,j=1}^p$ . The estimates of the population parameters are denoted as  $\hat{\Sigma}_n^{(k)} = (\hat{\sigma}_{ij}^{(k)})_{i,j=1}^p$ ,  $\hat{\Psi}_n^{(k)} = (\hat{\psi}_{n,ij}^{(k)})_{i,j=1}^p$ , and  $\hat{\Theta}_{\rho_n}^{(k)} = (\hat{\theta}_{\rho_n,ij}^{(k)})_{i,j=1}^p$ . For a vector  $x = (x_1, \dots, x_p)^T$  and  $J \subset \{1, \dots, p\}$ , we denote  $x_J = (x_j, j \in J)^T$ . For a matrix  $A$ ,  $\lambda_k(A)$  is the  $k$ th smallest eigenvalue and  $\vec{A}$  is the vectorization of  $A$ . For  $J \subset \{(i, j) : i, j = 1, \dots, p\}$  and  $A \in \mathbb{R}^{p \times p}$ ,  $\vec{A}_J$  is a vector in  $\mathbb{R}^{|J|}$  obtained by removing elements corresponding to  $(i, j) \notin J$  from  $\vec{A}$ . A zero-filled matrix  $A_J \in \mathbb{R}^{p \times p}$  is obtained from  $A$  by replacing  $a_{ij}$  by 0 for  $(i, j) \notin J$ .

### 8.1. Consistency in matrix norms

Theorem 1 is a direct consequence of the following result.

**Lemma 1.** (i) Suppose that Condition 1 holds. Let  $\gamma \in (0, \min_k \pi_k)$  be arbitrary. For

$$n \geq \max \left\{ \frac{6}{\gamma} \log p, \frac{2^{15} 3^3 C_1^2}{\gamma^3} (1 + 4c_1^2)^2 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2 \lambda_{\Theta}^4 \left(1 + \rho_2 \|L\|_2^{1/2}\right)^2 s \log p \right\}$$

and  $\rho_n = 2^3 \sqrt{6} C_1 (1 + 4c_1^2) \gamma^{-1/2} \max_{k,i} \sigma_{ii}^{(k)} \sqrt{\log p/n}$ , we have with probability  $(1 - 2K/p)(1 - 2K \exp(-2n(\min_k \pi_k - \gamma)^2))$  that

$$\sum_{k=1}^K \|\hat{\Theta}_{\rho_n}^{(k)} - \Theta_0^{(k)}\|_F \leq \frac{2^{15/2} 3^{3/2} C_1}{\gamma^{3/2}} (1 + 4c_1^2) \max_{k,i} \sigma_{ii}^{(k)} \lambda_{\Theta}^2 \left(1 + \rho_2 \|L\|_2^{1/2}\right) \sqrt{\frac{s \log p}{n}}.$$

(ii) Suppose that Condition 2 holds with  $p \leq c_7 n^{c_2}$ ,  $c_2, c_3, c_7 > 0$ . For  $\rho_n = C_1 K \delta_n$  satisfying

$$2^4 3^2 C_1 \rho_n^2 \gamma^{-2} s (1 + \rho_2 \|L\|_2^{1/2})^2 \lambda_{\Theta}^4 \leq 1/4$$

and  $\tau > (2^7 + 2^3 \sqrt{1 + 2^4 3^2 c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2}) / (9c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2)$  we have with probability  $(1 - 2K \exp(-2n(\min_k \pi_k - \gamma)^2))\nu_n$  that

$$\sum_{k=1}^K \|\hat{\Theta}_{\rho_n}^{(k)} - \Theta_0^{(k)}\|_F \leq 2^4 3^{3/2} C_1 \gamma^{-2} K \left(1 + \rho_2 \|L\|_2^{1/2}\right) \lambda_{\Theta}^2 s^{1/2} \delta_n,$$

where

$$\begin{aligned} \delta_n &\equiv \max_{k,i} \{\sigma_{ii}^{(k)}\}^2 c_4 (4 + \tau) \gamma^{-1} \frac{\log p}{n} \\ &\quad + (1 + 2 \max_{k,i} |\mu^{(k),i}|) \sqrt{\max_{k,i} \{\sigma_{ii}^{(k)}\}^2 c_4 (4 + \tau) \gamma^{-1} \frac{\log p}{n}} \\ &\quad + 2 \max_{k,i,j} \mathbb{E} |X^{(k),i} X^{(k),j}| I \left( |X^{(k),i} X^{(k),j}| \geq \sqrt{\frac{\gamma n}{\log p}} \right) \\ &\quad + 4 \left\{ \max_{k,i} \mathbb{E} |X^{(k),i}| I \left( |X^{(k),i}| \geq \sqrt{\frac{\gamma n}{\log p}} \right) \right\}^2 \\ &\quad + 2(1 + 2 \max_{k,i} |\mu^{(k),i}|) \max_{k,i} \mathbb{E} |X^{(k),i}| I \left( |X^{(k),i}| \geq \sqrt{\frac{\gamma n}{\log p}} \right) \\ &= O \left( \sqrt{\frac{\log p}{n}} \right), \end{aligned}$$

and

$$\begin{aligned} \nu_n &\equiv \frac{3c_7 c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2 (\log p)^{c_2+c_3+1}}{\gamma^{c_3} n^{c_3}} + \frac{c_7 c_4 \max_{k,i} \sigma_{ii}^{(k)} (\log p)^{2(c_2+c_3+1)}}{n^{c_2+c_3+1}} \\ &\quad + 8p^2 \exp \left( - \frac{\max_{k,i} \sigma_{ii}^{(k)} c_4 (4 + \tau) \log p}{2 \max_{k,i} \sigma_{ii}^{(k)} c_4 + \sqrt{\max_{k,i} \{\sigma_{ii}^{(k)}\}^2 c_4 (64 + 16\tau) / 3}} \right) \\ &= o(1). \end{aligned}$$

Our proofs adopt several tools from Negahban et al. [20]. Note however that our penalty does not penalize the diagonal elements, and is hence a seminorm; thus, their results do not apply to our case. We first introduce several notations. To treat multiple precision matrices in a unified way, our parameter space is defined to be the set  $\tilde{\mathbb{R}}^{(pK) \times (pK)}$  of  $(pK) \times (pK)$  symmetric block diagonal matrices, where the  $k$ th diagonal block is a  $p \times p$  matrix corresponding to the precision matrix of subpopulation  $k$ . We write  $A \in \tilde{\mathbb{R}}^{(pK) \times (pK)}$  for a  $K$ -tuple  $(A^{(k)})_{k=1}^K$  of diagonal blocks  $A^{(k)} \in \mathbb{R}^{p \times p}$ . Note that for  $A, B \in \tilde{\mathbb{R}}^{(pK) \times (pK)}$ ,  $\langle A, B \rangle_{pK} = \sum_{k=1}^K \langle A^{(k)}, B^{(k)} \rangle_p$  where  $\langle \cdot, \cdot \rangle_p$  is the trace inner product on  $\mathbb{R}^{p \times p}$ . In this parameter space, we evaluate the following map from  $\tilde{\mathbb{R}}^{(pK) \times (pK)}$  to  $\mathbb{R}$  given by

$$f(\Delta) = -\tilde{\ell}_n(\Theta_0 + \Delta) + \tilde{\ell}_n(\Theta_0) + \rho_n \{r(\Theta_0 + \Delta) - r(\Theta_0)\},$$

where  $r : \tilde{\mathbb{R}}^{(pK) \times (pK)} \mapsto \mathbb{R}$  is given by  $r(\Theta) = \|\Theta\|_1 + \rho_2 \|\Theta\|_L$ . This map provides information on the behavior of our criterion function in the neighborhood of  $\Theta_0$ . A similar map with a different penalty was studied in Rothman et al. [26]. A key observation is that  $f(0) = 0$  and  $f(\hat{\Delta}_n) \leq 0$  where  $\hat{\Delta}_n = \hat{\Theta}_{\rho_n} - \Theta_0$ .

The following lemma provides a non-asymptotic bound on the Frobenius norm of  $\Delta$  (see Lemma 4 in Negahban et al. [21] for a similar lemma in a different context). Let  $S = \cup_{k=1}^K S^{(k)}$  be the union of the supports of  $\Omega_0^{(k)}$ . Define a model subspace  $\mathcal{M} = \{\Omega \in \tilde{\mathbb{R}}^{(pK) \times (pK)} : \omega_{ij}^{(k)} = 0, (i, j) \notin S, k = 1, \dots, K\}$  and its orthocomplement  $\mathcal{M}^\perp = \{\Omega \in \tilde{\mathbb{R}}^{(pK) \times (pK)} : \omega_{ij}^{(k)} = 0, (i, j) \in S, k = 1, \dots, K\}$  under the trace inner product in  $\tilde{\mathbb{R}}^{(pK) \times (pK)}$ . For  $A = (a_{ij})_{i,j=1}^{pK} \in \tilde{\mathbb{R}}^{(pK) \times (pK)}$ , we write  $A = A_{\mathcal{M}} + A_{\mathcal{M}^\perp}$  where  $A_{\mathcal{M}}$  and  $A_{\mathcal{M}^\perp}$  are the projection of  $A$  into  $\mathcal{M}$  and  $\mathcal{M}^\perp$ , in the Frobenius norm, respectively. In other words, the  $(i, j)$ -element of  $A_{\mathcal{M}}$  is  $a_{ij}$  if  $(i, j) \in S$  and zero otherwise, and the  $(i, j)$ -element of  $A_{\mathcal{M}^\perp}$  is  $a_{ij}$  if  $(i, j) \notin S$  and zero otherwise. Note that  $\Theta_0 \in \mathcal{M}$ . Define the set  $\mathcal{C} = \{\Delta \in \tilde{\mathbb{R}}^{(pK) \times (pK)} : r(\Delta_{\mathcal{M}^\perp}) \leq 3r(\Delta_{\mathcal{M}})\}$ .

**Lemma 2.** *Let  $\epsilon > 0$  be arbitrary. Suppose  $\rho_n \geq 2 \max_{1 \leq k \leq K} \|\hat{\Psi}_n^{(k)} - \Psi_0^{(k)}\|_\infty$ . If  $f(\Delta) > 0$  for all elements  $\Delta \in \mathcal{C} \cap \{\Delta \in \tilde{\mathbb{R}}^{(pK) \times (pK)} : \|\Delta\|_F = \epsilon\}$  then  $\|\hat{\Delta}_n\|_F \leq \epsilon$ .*

*Proof.* We first show that  $\hat{\Delta}_n \in \mathcal{C}$ . We have by the convexity of  $-\tilde{\ell}_n(\Theta)$  that

$$-\tilde{\ell}_n(\Theta_0 + \hat{\Delta}_n) + \tilde{\ell}_n(\Theta_0) \geq -|\langle -\nabla \tilde{\ell}_n(\Theta_0), \hat{\Delta}_n \rangle|.$$

It follows from Lemma 3(iv) with our choice  $\rho_n$  that the right hand side of the inequality is further bounded below by  $-2^{-1}\rho_n \left( r(\hat{\Delta}_{n,\mathcal{M}}) + r(\hat{\Delta}_{n,\mathcal{M}^\perp}) \right)$ . Applying Lemma 3(iii), we obtain

$$\begin{aligned} 0 &\geq f(\hat{\Delta}_n) = -\tilde{\ell}_n(\Theta_0 + \hat{\Delta}_n) + \tilde{\ell}_n(\Theta_0) + r(\Theta_0 + \hat{\Delta}_n) - r(\Theta_0) \\ &\geq \frac{\rho_n}{2} r(\hat{\Delta}_{n,\mathcal{M}^\perp}) - \frac{3\rho_n}{2} r(\hat{\Delta}_{n,\mathcal{M}}), \end{aligned}$$

or  $r(\hat{\Delta}_{n,\mathcal{M}^\perp}) \leq 3r(\hat{\Delta}_{n,\mathcal{M}})$ . This verifies  $\hat{\Delta}_n \in \mathcal{C}$ . Note that  $f$ , as a function of  $\Delta$  is sum of two convex functions  $\ell_n$  and  $r$ , and is hence convex. Thus, the rest of the proof follows exactly as Lemma 4 in Negahban et al. [21].  $\square$

**Lemma 3.** *Let  $\Delta \in \tilde{\mathbb{R}}^{(pK) \times (pK)}$ .*

(i) *The gradient of  $\tilde{\ell}_n(\Theta_0)$  is a block diagonal matrix given by*

$$\nabla \tilde{\ell}_n(\Theta_0) = n^{-1} \text{diag}\{n_1(\Psi_0^{(1)} - \hat{\Psi}_n^{(1)}), \dots, n_K(\Psi_0^{(K)} - \hat{\Psi}_n^{(K)})\}. \quad (10)$$

(ii) *Let  $c > 0$  be a constant. For  $\|\Delta\|_F \leq c$  and  $n_k/n \geq \gamma > 0$  for all  $k$  and  $n$ ,*

$$-\tilde{\ell}_n(\Theta_0 + \Delta) + \tilde{\ell}_n(\Theta_0) + \langle \nabla \tilde{\ell}_n(\Theta_0), \Delta \rangle \geq \frac{\gamma}{2\{\lambda_\Theta + c\}^2} \|\Delta\|_F^2 \equiv \kappa_{\ell_n, c} \|\Delta\|_F^2. \quad (11)$$

(iii) *The map  $r$  is a seminorm, convex, and decomposable with respect to  $(\mathcal{M}, \mathcal{M}^\perp)$  in the sense that  $r(\Theta_1 + \Theta_2) = r(\Theta_1) + r(\Theta_2)$  for every  $\Theta_1 \in \mathcal{M}$*

and  $\Theta_2 \in \mathcal{M}^\perp$ . Moreover,

$$r(\Theta_0 + \Delta) - r(\Theta_0) \geq r(\Delta_{\mathcal{M}^\perp}) - r(\Delta_{\mathcal{M}}).$$

(iv) For  $\Delta \in \tilde{\mathbb{R}}^{(pK) \times (pK)}$ ,

$$|\langle \nabla \tilde{\ell}_n(\Theta_0), \Delta \rangle| \leq r(\Delta) \max_{1 \leq k \leq K} \|\hat{\Psi}_n^{(k)} - \Psi_0^{(k)}\|_\infty. \quad (12)$$

(v) For  $\Theta \in \tilde{\mathbb{R}}^{(pK) \times (pK)}$ ,

$$r(\Theta_{\mathcal{M}}) \leq (s+1)^{1/2} \left(1 + \rho_2 \|L\|_2^{1/2}\right) \|\Theta_{\mathcal{M}}\|_F.$$

*Proof.* (i) The result follows by taking derivatives blockwise.

(ii) Rothman et al. [26] (page 500-502) showed that

$$\begin{aligned} & -\tilde{\ell}_n(\Theta_0 + \Delta) + \tilde{\ell}_n(\Theta_0) - \langle -\nabla \tilde{\ell}_n(\Theta_0), \Delta \rangle \\ &= \sum_{k=1}^K \frac{n_k}{n} \left( -\log \det(\Theta_0^{(k)} + \Delta^{(k)}) + \log \det(\Theta_0^{(k)}) + \langle \Psi_0^{(k)}, \Delta^{(k)} \rangle \right) \\ &\geq \sum_{k=1}^K \frac{n_k}{n} \frac{\|\Delta^{(k)}\|_F^2}{2 \min_{0 \leq v \leq 1} \left\{ \left\| \Theta_0^{(k)} \right\|_2 + v \left\| \Delta^{(k)} \right\|_2 \right\}^2}. \end{aligned}$$

Since  $\|A\|_2 \leq \|A\|_F$ ,  $n_k/n \geq \gamma$  and  $\|\Delta\|_F \leq c$ , this is further bounded below by

$$\sum_{k=1}^K \frac{\gamma}{2} \frac{\|\Delta^{(k)}\|_F^2}{\left\{ \left\| \Theta_0^{(k)} \right\|_2 + \left\| \Delta^{(k)} \right\|_F \right\}^2} \geq \kappa_{\ell_n, c} \|\Delta\|_F^2.$$

(iii) Because the graph Laplacian  $L$  is a positive semidefinite matrix, the triangle inequality  $r(\Theta_1 + \Theta_2) \leq r(\Theta_1) + r(\Theta_2)$  holds. To see this let  $L = \tilde{L}\tilde{L}^T$  be any Cholesky decomposition of  $L$ . Then

$$\{(x+y)^T L(x+y)\}^{1/2} = \|\tilde{L}^T(x+y)\| \leq \|\tilde{L}^T x\| + \|\tilde{L}^T y\| = \{x^T L x\}^{1/2} + \{y^T L y\}^{1/2}.$$

It is clear that  $r(c\Theta) = cr(\Theta)$  for any constant  $c$ . Thus, given that  $r$  does not penalize the diagonal elements, it is a seminorm. The decomposability follows from the definition of  $r$ . The convexity follows from the same argument for the triangle inequality. Since  $\Theta_0 + \Delta = \Theta_0 + \Delta_{\mathcal{M}} + \Delta_{\mathcal{M}^\perp}$ , the triangle inequality and the decomposability of  $r$  yield

$$r(\Theta_0 + \Delta) - r(\Theta_0) \geq r(\Theta_0 + \Delta_{\mathcal{M}^\perp}) - r(\Delta_{\mathcal{M}}) - r(\Theta_0) = r(\Delta_{\mathcal{M}^\perp}) - r(\Delta_{\mathcal{M}}).$$

(iv) We show that, for  $A, B \in \tilde{\mathbb{R}}^{(pK) \times (pK)}$  with  $\text{diag}(B) = 0$ ,  $\langle A, B \rangle \leq r(A)\|B\|_\infty$ . If  $A$  is a diagonal matrix (or if  $A = 0$ ), the inequality trivially holds since  $\langle A, B \rangle = 0$ . If not,  $r(A) \neq 0$  so that

$$\frac{\langle A, B \rangle}{r(A)} \leq \frac{\|A\|_1 \|B\|_\infty}{\|A\|_1} = \|B\|_\infty.$$

Since the diagonal elements of  $\nabla \tilde{\ell}_n(\Theta_0)$  are all zero, the result follows.

(v) For  $s \neq 0$ , we have

$$\begin{aligned} \frac{r(\Theta_{\mathcal{M}})}{\|\Theta_{\mathcal{M}}\|_F} &\leq \sup_{\Theta \in \mathcal{M}} \frac{\sum_{k=1}^K \|\Theta^{(k)}\|_1}{\|\Theta - \text{diag}(\Theta)\|_F} + \sup_{\Theta \in \mathcal{M}} \frac{\rho_2 \sum_{i \neq j} \sqrt{\theta_{ij}^T L \theta_{ij}}}{\|\Theta\|_F} \\ &\leq s^{1/2} + \rho_2 \sup_{\Theta \in \mathcal{M}} \frac{\sum_{i \neq j} \sqrt{\|L\|_2 \|\theta_{ij}\|_F^2}}{\|\Theta\|_F} \\ &\leq s^{1/2} \left(1 + \rho_2 \|L\|_2^{1/2}\right). \end{aligned}$$

In the last inequality we used that  $\sqrt{\sum_{j=1}^J \sum_{i=1}^I a_{ij}^2} \geq J^{-1/2} \sum_{j=1}^J \sqrt{\sum_{i=1}^I a_{ij}^2}$ , which follows by the concavity of the square root function. For  $s = 0$ , we trivially have  $0 = r(\Theta_{\mathcal{M}}) \leq s^{1/2} \{1 + \rho_2 \|L\|_2^{1/2}\} \|\Theta_{\mathcal{M}}\|_F$ . Combining these two cases yields the desired result.  $\square$

Next, we obtain an upper bound for  $\max_{1 \leq k \leq K} \|\hat{\Psi}_n^{(k)} - \Psi_0^{(k)}\|_\infty$ , which holds with high-probability assuming the tail conditions of the random vectors.

**Lemma 4.** Suppose that  $n_k/n \geq \gamma > 0$  for all  $k$  and  $n$ .

(i) Suppose that Condition 1 holds. Then for  $n \geq 6\gamma^{-1} \log p$  we have

$$P\left(\|\hat{\Sigma}_n - \Sigma_0\|_\infty \geq 2^3 \sqrt{6} (1 + 4c_1^2)^2 \gamma^{-1/2} \max_{k,i} \sigma_{ii}^{(k)} \sqrt{\frac{\log p}{\gamma n}}\right) \leq 2K/p. \quad (13)$$

(ii) Suppose that Condition 2 holds with  $c_2, c_3 > 0$  and  $p \leq c_7 n^{c_2}$ . Then we have for  $\tau > \max_k (2^7 + 2^3 \sqrt{1 + 2^4 3^2 c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2}) / (9c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2)$

$$P\left(\|\hat{\Sigma}_n - \Sigma_0\|_\infty \geq \sum_{k=1}^K \delta_n^{(k)}\right) \leq K\nu_n \quad (14)$$

where

$$\delta_n^{(k)} \equiv (1 + 2 \max_i |\mu^{(k),i}|) (2\delta_{n,1}^{(k)} + \delta_{n,2}^{(k)}) + (\delta_{n,1}^{(k)})^2 + (\delta_{n,2}^{(k)})^2 + 2\delta_{n,3}^{(k)},$$

with

$$\delta_{n,1}^{(k)} \equiv \max_{i,j} \mathbb{E} |X_l^{(k),i} X_l^{(k),j}| I(|X_l^{(k),i} X_l^{(k),j}| \geq n_k^{1/2} (\log p)^{-1/2}),$$

$$\delta_{n,2}^{(k)} \equiv \{c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}\}^2 (4 + \tau) \log p / n_k^{1/2},$$

$$\delta_{n,3}^{(k)} \equiv \max_i \mathbb{E} |X_l^{(k),i}| I(|X_l^{(k),i}| \geq n_k^{1/2} (\log p)^{-1/2}).$$

(iii) Suppose that Condition 3 holds and that  $P(\|\hat{\Sigma}_n - \Sigma_0\|_\infty \geq b_n) = o(1)$  and  $b_n = o(1)$  as  $n \rightarrow \infty$ . Then  $P(\|\hat{\Psi}_n - \Psi_0\|_\infty \geq C_1 b_n) = o(1)$ .

*Proof.* (i) This was proved by Ravikumar et al. [25].

(ii) Note that

$$\begin{aligned}\hat{\Sigma}_n^{(k)} - \Sigma^{(k)} &= n_k^{-1} \sum_{l=1}^{n_k} (X_l^{(k)})^{\otimes 2} - \mathbb{E}(X^{(k)})^{\otimes 2} - (\bar{X}^{(k)} - \mu^{(k)})^{\otimes 2} \\ &\quad - \mu^{(k)}(\bar{X}^{(k)} - \mu^{(k)})^T - (\bar{X}^{(k)} - \mu^{(k)})(\mu^{(k)})^T.\end{aligned}$$

We first evaluate the probability in (14) for  $n_k^{-1} \sum_{l=1}^{n_k} (X_l^{(k)})^{\otimes 2} - \mathbb{E}(X^{(k)})^{\otimes 2}$ . Let

$$\begin{aligned}Y_l^{(k),ij} &\equiv X_l^{(k),i} X_l^{(k),j} - \mathbb{E} X_l^{(k),i} X_l^{(k),j}, \\ \bar{Y}_l^{(k),ij} &\equiv X_l^{(k),i} X_l^{(k),j} I \left( |X_l^{(k),i} X_l^{(k),j}| \leq \sqrt{\frac{n_k}{\log p}} \right) \\ &\quad - \mathbb{E} X_l^{(k),i} X_l^{(k),j} I \left( |X_l^{(k),i} X_l^{(k),j}| \leq \sqrt{\frac{n_k}{\log p}} \right), \\ \tilde{Y}_l^{(k),ij} &\equiv Y_l^{(k),ij} - \bar{Y}_l^{(k),ij}.\end{aligned}$$

We have

$$\begin{aligned}&P \left( \max_{i,j} \left| \sum_{l=1}^{n_k} \tilde{Y}_l^{(k),ij} \right| \geq 2n_k \delta_{n,1}^{(k)} \right) \\ &\leq P \left( \max_{i,j} \left| \sum_{l=1}^{n_k} X_l^{(k),i} X_l^{(k),j} I \left( |X_l^{(k),i} X_l^{(k),j}| \geq \sqrt{\frac{n_k}{\log p}} \right) \right| \geq n_k \delta_{n,1}^{(k)} \right) \\ &\leq P \left( \max_{l,i} (X_l^{(k),i})^2 \geq n_k^{1/2} (\log p)^{-1/2} \right) \quad (xy \leq \max\{x^2, y^2\}) \\ &\leq pn_k \frac{\mathbb{E} X_{0i}^{4(c_2+c_3+1)} (\log p)^{c_2+c_3+1}}{n_k^{c_2+c_3+1}} \quad (\text{Markov's inequality}) \\ &\leq \frac{c_7 c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2 (\log p)^{c_2+c_3+1}}{n_k^{c_3}} \quad (p \leq c_7 n^{c_2}) \\ &\leq \frac{c_7 c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2 (\log p)^{c_2+c_3+1}}{\gamma^{c_3} n^{c_3}} \equiv \nu_{n,1},\end{aligned} \tag{15}$$

where the first inequality follows from the triangle inequality. Note that

$$\begin{aligned}\mathbb{E} \left( \bar{Y}_l^{(k),ij} \right)^2 &\leq \mathbb{E} \left[ X_l^{(k),i} X_l^{(k),j} I \left( |X_l^{(k),i} X_l^{(k),j}| \leq \sqrt{\frac{n_k}{\log p}} \right) \right]^2 \leq \mathbb{E} |X_l^{(k),i} X_l^{(k),j}|^2 \\ &\leq 2^{-1} (\mathbb{E} (X_l^{(k),i})^4 + \mathbb{E} (X_l^{(k),j})^4) \leq c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2.\end{aligned}$$

It follows from Bernstein's inequality that

$$P \left( \max_{i,j} \left| \sum_{l=1}^{n_k} \bar{Y}_l^{(k),ij} \right| \geq n_k \delta_{n,2}^{(k)} \right)$$

$$\leq 2p^2 \exp \left( - \frac{c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2 (4 + \tau) \log p}{2c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2 + 2\sqrt{c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2 (64 + 16\tau)/3}} \right) \equiv \nu_{n,2}. \quad (16)$$

Now, for  $\tau > (2^7 + 2^3 \sqrt{1 + 2^4 3^2 c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2}) / (9c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2)$ ,  $\nu_{n,2} \rightarrow 0$  as  $p \rightarrow \infty$ . Note that for this to hold it suffices to have

$$\frac{3c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2 (4 + \tau)}{6c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2 + 8\sqrt{c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2 (4 + \tau)}} > 2,$$

so that the power in the exponent is negative. This inequality reduces to

$$3c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2 \tau > 16\sqrt{c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2 (4 + \tau)}.$$

We can solve this by changing a quadratic equation for  $\tau$ , since  $\tau$  of our interest is positive. Combining (15) and (16) yields

$$P \left( \left\| \frac{1}{n_k} \sum_{i=1}^{n_k} (X_l^{(k)})^{\otimes 2} - \mathbb{E}(X^{(k)})^{\otimes 2} \right\|_{\infty} \geq 2\delta_{n,1}^{(k)} + \delta_{n,2}^{(k)} \right) \leq \nu_{n,1} + \nu_{n,2}. \quad (17)$$

Let

$$Z_l^{(k),i} \equiv X_l^{(k),i} - \mathbb{E} X_l^{(k),i},$$

$$\bar{Z}_l^{(k),i} \equiv X_l^{(k),i} I(|X_l^{(k),i}| \leq n_k^{1/2} (\log p)^{-1/2}) - \mathbb{E} X_l^{(k),i} I(|X_l^{(k),i}| \leq n_k^{1/2} (\log p)^{-1/2}),$$

$$\tilde{Z}_l^{(k),i} \equiv U_l^{(k),i} - \bar{Z}_l^{(k),i}.$$

Proceeding as for  $Y_l^{(k),ij}$ 's, we have

$$P \left( \max_i \left| \sum_{l=1}^{n_k} \tilde{Z}_l^{(k),i} \right| \geq 2n_k \delta_{n,3}^{(k)} \right) \leq \frac{c_7 c_4 \max_{k,i} \{\sigma_{ii}^{(k)}\}^2 (\log p)^{2(c_2 + c_3 + 1)}}{\gamma^{c_2 + c_3 + 1} n^{c_2 + c_3 + 1}} \equiv \nu_{n,3},$$

and

$$P \left( \max_i \left| \sum_{k=1}^n \bar{Z}_l^{(k),i} \right| \geq n_k \delta_{n,2}^{(k)} \right) \leq \nu_{n,2}.$$

Thus, we have

$$\begin{aligned} P(\|\bar{X}^{(k)} - \mu^{(k)}\|_{\infty}^{\otimes 2} \geq (\delta_{n,2}^{(k)})^2 + (2\delta_{n,3}^{(k)})^2) \\ \leq P \left( \max_i |\bar{X}^{(k),i} - \mu^{(k),i}| \geq \sqrt{(\delta_{n,1}^{(k)})^2 + (\delta_{n,2}^{(k)})^2} \right) \\ \leq P \left( \max_i \left| \sum_{k=1}^n \bar{Z}_l^{(k),i} \right| \geq n_k \delta_{n,2}^{(k)} \right) + P \left( \max_i \left| \sum_{l=1}^{n_k} \tilde{Z}_l^{(k),i} \right| \geq 2n_k \delta_{n,3}^{(k)} \right) \\ \leq \nu_{n,2} + \nu_{n,3}, \end{aligned} \quad (18)$$

and

$$\begin{aligned} P\left(\|(\bar{X}^{(k)} - \mu^{(k)})(\mu^{(k)})^T\|_\infty \geq \max_i |\mu^{(k),i}|(2\delta_{n,1}^{(k)} + \delta_{n,2}^{(k)})\right) \\ \leq P\left(\max_i |\bar{X}^{(k),i} - \mu^{(k),i}| \geq 2\delta_{n,1}^{(k)} + \delta_{n,2}^{(k)}\right) \leq \nu_{n,1} + \nu_{n,2}. \end{aligned} \quad (19)$$

Combining (17)–(19) yields

$$\begin{aligned} P\left(\|\hat{\Sigma}_n^{(k)} - \Sigma^{(k)}\|_\infty \geq (1 + 2\max_i |\mu^{(k),i}|)(2\delta_{n,1}^{(k)} + \delta_{n,2}^{(k)}) + (\delta_{n,2}^{(k)})^2 + (2\delta_{n,3}^{(k)})^2\right) \\ \leq 3\nu_{n,1} + 4\nu_{n,2} + \nu_{n,3} = \nu_n. \end{aligned}$$

Note that  $\delta_{n,1}^{(k)}, \delta_{n,2}^{(k)}, \delta_{n,3}^{(k)}, \nu_{n,1}, \nu_{n,2}, \nu_{n,3} \rightarrow 0$  as  $n, p \rightarrow \infty$  if  $\log p/n \rightarrow 0$ . Note also that  $\delta_{n,1}^{(k)}, \delta_{n,2}^{(k)}$  and  $(\delta_{n,3}^{(k)})^2$  are  $O(\sqrt{\log p/n})$  on the set where  $n_k/n \geq \gamma$ . For example, we have by Jensen's inequality that

$$\begin{aligned} \sqrt{\frac{n}{\log p}}(\delta_{n,3}^{(k)})^2 &= \sqrt{\frac{n}{\log p}} \max_i \{\mathbb{E} |X^{(k),i}| I\{|X^{(k),i}| \geq n_k^{1/2}(\log p)^{-1/2}\}\}^2 \\ &\leq \max_i \mathbb{E} \frac{n}{n_k} \sqrt{\frac{n_k}{\log p}} |X^{(k),i}|^2 I\{|X^{(k),i}| \geq n_k^{1/2}(\log p)^{-1/2}\} \\ &\leq \gamma^{-1} \max_i \mathbb{E} |X^{(k),i}|^3 I\{|X^{(k),i}| \geq n_k^{1/2}(\log p)^{-1/2}\} \\ &\leq c_4 \gamma^{-1} \max_i \{\sigma_{ii}^{(k)}\}^2. \end{aligned}$$

(iii) Given that  $|\sigma_{0,ij}^{(k)}| \leq \sqrt{\sigma_{0,ii}^{(k)}\sigma_{0,jj}^{(k)}}$ ,

$$\begin{aligned} |\psi_{n,ij}^{(k)} - \psi_{0,ij}^{(k)}| &= \left| \frac{\hat{\sigma}_{n,ij}^{(k)}}{\sqrt{\hat{\sigma}_{n,ii}^{(k)}\hat{\sigma}_{n,jj}^{(k)}}} - \frac{\sigma_{0,ij}^{(k)}}{\sqrt{\sigma_{0,ii}^{(k)}\sigma_{0,jj}^{(k)}}} \right| \\ &= \frac{\left| \sqrt{\sigma_{0,ii}^{(k)}\sigma_{0,jj}^{(k)}}(\hat{\sigma}_{n,ij}^{(k)} - \sigma_{0,ij}^{(k)}) + \sigma_{0,ij}^{(k)} \left( \sqrt{\sigma_{0,ii}^{(k)}\sigma_{0,jj}^{(k)}} - \sqrt{\hat{\sigma}_{n,ii}^{(k)}\hat{\sigma}_{n,jj}^{(k)}} \right) \right|}{\sqrt{\hat{\sigma}_{n,ii}^{(k)}\hat{\sigma}_{n,jj}^{(k)}\sigma_{0,ii}^{(k)}\sigma_{0,jj}^{(k)}}} \\ &\leq \frac{\sqrt{\sigma_{0,ii}^{(k)}\sigma_{0,jj}^{(k)}}}{\sqrt{\hat{\sigma}_{n,ii}^{(k)}\hat{\sigma}_{n,jj}^{(k)}\sigma_{0,ii}^{(k)}\sigma_{0,jj}^{(k)}}} \left\{ \left| \hat{\sigma}_{n,ij}^{(k)} - \sigma_{0,ij}^{(k)} \right| + \left| \sqrt{\sigma_{0,ii}^{(k)}\sigma_{0,jj}^{(k)}} - \sqrt{\hat{\sigma}_{n,ii}^{(k)}\hat{\sigma}_{n,jj}^{(k)}} \right| \right\}, \end{aligned}$$

wherein

$$\begin{aligned} &\sqrt{\sigma_{0,ii}^{(k)}\sigma_{0,jj}^{(k)}} - \sqrt{\hat{\sigma}_{n,ii}^{(k)}\hat{\sigma}_{n,jj}^{(k)}} \\ &= \frac{\sqrt{\sigma_{0,jj}^{(k)}}}{\sqrt{\sigma_{0,ii}^{(k)}} + \sqrt{\hat{\sigma}_{n,ii}^{(k)}}} (\sigma_{0,ii}^{(k)} - \hat{\sigma}_{n,ii}^{(k)}) + \frac{\sqrt{\hat{\sigma}_{n,ii}^{(k)}}}{\sqrt{\sigma_{0,jj}^{(k)}} + \sqrt{\hat{\sigma}_{n,jj}^{(k)}}} (\sigma_{0,jj}^{(k)} - \hat{\sigma}_{n,jj}^{(k)}). \end{aligned}$$

Since  $b_n \rightarrow 0$ ,  $b_n \leq c_5/2$  for  $n$  sufficiently large by Condition 3. On the event  $\|\hat{\Sigma}_n - \Sigma_0\|_\infty \leq b_n$  with  $n$  large,  $0 < c_5/2 \leq \sigma_{0,ii}^{(k)} - c_5/2 \leq \hat{\sigma}_{n,ii}^{(k)} \leq \sigma_{0,ii}^{(k)} + c_5/2 \leq$



$c_6 + c_5/2$ . Thus,

$$\begin{aligned} \frac{\sqrt{\sigma_{0,ii}^{(k)}\sigma_{0,jj}^{(k)}}}{\sqrt{\hat{\sigma}_{n,ii}^{(k)}\hat{\sigma}_{n,jj}^{(k)}\sigma_{0,ii}^{(k)}\sigma_{0,jj}^{(k)}}} &\leq \frac{2(c_5 + 2c_6)}{c_5^2} \\ \frac{\sqrt{\sigma_{0,jj}^{(k)}}}{\sqrt{\sigma_{0,ii}^{(k)} + \sqrt{\hat{\sigma}_{n,ii}^{(k)}}}} &\leq \frac{\sqrt{c_6}}{2\sqrt{c_5}} \\ \frac{\sqrt{\hat{\sigma}_{n,ii}^{(k)}}}{\sqrt{\sigma_{0,jj}^{(k)} + \sqrt{\hat{\sigma}_{n,jj}^{(k)}}}} &\leq \frac{\sqrt{c_5 + 2c_6}}{2\sqrt{c_5}}. \end{aligned}$$

It follows that

$$\begin{aligned} |\psi_{n,ij}^{(k)} - \psi_{0,ij}^{(k)}| &\leq \\ &\left\{ 2c_5^{-2} + c_5 + c_6^{-3/2} + 2c_5^{-5/2}c_6 + (c_5^{-4} + 2c_5^{-5}c_6)^{1/2} \right\} \max_{k,i,j} |\hat{\sigma}_{n,ij}^{(k)} - \sigma_{0,ij}^{(k)}|. \end{aligned}$$

Thus, we have

$$\begin{aligned} P\left(\|\hat{\Psi}_n - \Psi_0\|_\infty \geq C_1 b_n\right) &\leq P\left(\|\hat{\Psi}_n - \Psi_0\|_\infty \geq C_1 b_n, \|\hat{\Sigma}_n - \Sigma_0\|_\infty < b_n\right) + P\left(\|\hat{\Sigma}_n - \Sigma_0\|_\infty \geq b_n\right) \\ &\leq 2P\left(\|\hat{\Sigma}_n - \Sigma_0\|_\infty \geq b_n\right) \rightarrow 0. \quad \square \end{aligned}$$

So far we have assumed  $n_k/n \geq \gamma$  in lemmas. We evaluate the probability of this event noting that  $n_k \sim \text{Binom}(n, \pi_k)$ .

**Lemma 5.** *Let  $\epsilon > 0$  such that  $\gamma \equiv \min_k \pi_k - \epsilon > 0$ . Then*

$$P\left(\min_k n_k/n \leq \min_k \pi_k - \epsilon\right) \leq 2K \exp(-2n\epsilon^2). \quad (20)$$

*Proof.* We have by Hoeffding's inequality that

$$\begin{aligned} P\left(\min_k n_k/n \leq \min_k \pi_k - \epsilon\right) &\leq P\left(\exists k, n_k/n \leq \min_k \pi_k - \epsilon\right) \\ &\leq P(\exists k, n_k/n \leq \pi_k - \epsilon) \leq P(\exists k, |n_k/n - \pi_k| \geq \epsilon) \\ &\leq \sum_{k=1}^K P(|n_k/n - \pi_k| \geq \epsilon) \leq 2K \exp(-2n\epsilon^2). \quad \square \end{aligned}$$

*Proof of Lemma 1.* We apply Lemma 2 to obtain the non-asymptotic error bounds.

We first compute a lower bound for  $f(\Delta)$ . Suppose  $\epsilon \leq c$ . For  $\Delta \in \mathcal{C} \cap \{\Delta \in \tilde{\mathbb{R}}^{(pK) \times (pK)} : \|\Delta\|_F = \epsilon\}$ , we have by Lemma 3(ii) and (iii) that

$$f(\Delta) \geq -\langle \tilde{\ell}_n(\Theta_0), \Delta \rangle + \kappa_{\ell_n, c} \|\Delta\|_F^2 + \rho_n \{r(\Delta_{\mathcal{M}^\perp}) - r(\Delta_{\mathcal{M}})\}.$$

The assumption on  $\rho_n$  and Lemma 3(iii) and (iv) then yield

$$|\langle \tilde{\ell}_n(\Theta_0), \Delta \rangle| \leq \frac{\rho_n}{2} \{r(\Delta_{\mathcal{M}}) + r(\Delta_{\mathcal{M}^\perp})\}.$$

From this inequality and Lemma 3(v) we have

$$\begin{aligned} f(\Delta) &\geq \kappa_{\ell_n, c} \|\Delta\|_F^2 - \frac{3\rho_n}{2} r(\Delta_{\mathcal{M}}) \\ &\geq \kappa_{\ell_n, c} \|\Delta\|_F^2 - \frac{3\rho_n}{2} (s+1)^{1/2} \left(1 + \rho_2 \|L\|_2^{1/2}\right) \|\Delta\|_F. \end{aligned}$$

Viewing the right hand side of the above inequality as a quadratic equation in  $\|\Delta\|_F$ , we have  $f(\Delta) > 0$  if

$$\|\Delta\|_F \geq \frac{3\rho_n}{\kappa_{\ell_n, c}} (s+1)^{1/2} \left(1 + \rho_2 \|L\|_2^{1/2}\right) \equiv \epsilon_c > 0.$$

Thus, if we show that there exists a  $c_0 > 0$  such that  $\epsilon_{c_0} \leq c_0$ , Lemma 2 yields that  $\|\hat{\Theta}_{\rho_n} - \Theta_0\|_F \leq \epsilon_{c_0}$ .

Consider the inequality  $(x+y)^2 z^{1/2} \leq y$  where  $x, y, z \geq 0$ . This inequality holds for  $(x, y, z)$  such that  $x = y$  and  $xz^{1/2} = 1/4$ . We apply the inequality above with  $x = \lambda_\Theta, y = c, z = 2^4 3^2 \rho_n^2 \gamma^{-2} s (1 + \rho_2 \|L\|_2^{1/2})^2$  and solve  $xz \leq 1/4$  for  $n$ . (i) For  $\rho_n = 2^3 \sqrt{6} C_1 (1 + 4c_1^2)^2 \gamma^{-1/2} \max_{k,i} \sigma_{ii}^{(k)} \sqrt{\log p/n}$ ,  $xz \leq 1/4$  yields

$$n \geq \max \left\{ \frac{6 \log p}{\gamma}, \frac{2^{15} 3^3 C_1^2 (1 + 4c_1^2)^2}{\gamma^3} \max_{k,i} \{\sigma_{ii}^{(k)}\}^2 \lambda_\Theta^4 \left(1 + \rho_2 \|L\|_2^{1/2}\right)^2 s \log p \right\},$$

and  $(x+y)^4 z$  becomes

$$\epsilon_{\max_k \{\|\Theta_0^{(k)}\|_2\}}^2 \leq 2^{15} 3^3 (1 + c_1^2)^2 \max_{k,i} (\sigma_{ii}^{(k)})^2 \left(1 + \rho_2 \|L\|_2^{1/2}\right)^2 \gamma^{-3} \lambda_\Theta^4 \frac{s \log p}{n}.$$

(ii) For  $\rho_n = C_1 K \delta_n$ , there is no closed form solution for  $n$ . Note that  $\delta_n \rightarrow 0$  if  $\log p/n \rightarrow 0$  so that  $xz \leq 1/4$  holds for  $n$  sufficiently large, given that  $\sum_{k=1}^K \delta_n^{(k)} \leq K \delta_n$ .

Computing appropriate probabilities using Lemmas 4 and 5 completes the proof.  $\square$

*Proof of Theorem 1.* The estimation error  $\|\hat{\Omega}_{\rho_n}^{(k)} - \Omega_0\|_2^{(2)}$  in the spectral norm can be bounded and evaluated in the same way as in the proof of Theorem 2 of Rothman et al. [26] together with Lemma 1.  $\square$

## 8.2. Model selection consistency

Our proof is based on the primal-dual witness approach of Ravikumar et al. [25], with some modifications to overcome a difficulty in their proof when applying

the fixed point theorem to a discontinuous function. First, we define the oracle estimator  $\check{\Theta}_{\rho_n} = (\check{\Theta}_{\rho_n}^{(1)}, \dots, \check{\Theta}_{\rho_n}^{(K)})$  by

$$\begin{aligned} \check{\Theta}_{\rho_n} = & \arg \min_{\Theta^{(k)} > 0, \Theta^{(k)} = (\Theta^{(k)})^T, \Theta_{(S^{(k)})^c}^{(k)} = 0} n^{-1} \sum_{k=1}^K n_k \left( \text{tr} \left( \Psi_n^{(k)} \Theta^{(k)} \right) - \log \det(\Theta^{(k)}) \right) \\ & + \rho_n \sum_{k=1}^K \|\Theta^{(k)}\|_1 + \rho_n \rho_2 \sum_{i,j} \sqrt{\Theta_{ij}^T L \Theta_{ij}}, \end{aligned} \quad (21)$$

where  $\Theta_{(S^{(k)})^c}^{(k)} = 0$  indicates that  $\Theta_{(i,j)}^{(k)} = 0$  for  $(i, j) \notin S^{(k)}$ .

**Lemma 6.** (i) Let  $A \in \mathbb{R}^{p \times p}$  be a positive semidefinite matrix with eigenvalues  $0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_p$  and corresponding eigenvectors  $u_i$  satisfying  $u_i \perp u_j, i \neq j$  and  $\|u_i\| = 1$ . The subdifferential  $\partial \sqrt{x^T A x}$  of  $f(x) = \sqrt{x^T A x}$  is

$$\partial \sqrt{x^T A x} = \begin{cases} Ax / \sqrt{x^T A x}, & Ax \neq 0, \\ \{U \Lambda^{1/2} y : \|y\| \leq 1\}, & Ax = 0. \end{cases}$$

where  $U \in \mathbb{R}^{p \times p}$  has  $u_i$  as the  $i$ th columns and  $\Lambda^{1/2}$  is the diagonal matrix with  $\lambda_i^{1/2}, i = 1, \dots, p$ , as diagonal elements. Furthermore, the subgradients are bounded above, i.e.

$$\|\nabla f(x)\|_\infty \leq \|A\|_2^{1/2}, \quad \text{for all } \nabla f(x) \in \partial \sqrt{x^T A x}.$$

(ii) Let  $A \in \mathbb{R}^{p \times p}$  be a positive semidefinite matrix and  $S = \{S_i\} \subset \{1, \dots, p\}$ . Suppose  $A_{SS}$  has eigenvalues  $0 \leq \lambda_{1,S} \leq \lambda_{2,S} \leq \dots \leq \lambda_{|S|,S}$  and corresponding eigenvectors  $u_{i,S}$  satisfying  $u_{i,S} \perp u_{j,S}, i \neq j$  and  $\|u_{i,S}\| = 1$ . Let  $g_S : \mathbb{R}^{|S|} \rightarrow \mathbb{R}^p$  be a map defined by  $g_S(x) = y$  where  $y_i = x_{S_j}$  for  $i = S_j$  for and  $y_i = 0$  for  $i \notin S$ . The subdifferential  $h_{A,S}(x) = \sqrt{g_S(x)^T A g_S(x)}$  equals to the subdifferential of  $\sqrt{x^T A_{SS} x}$  given by

$$\partial \sqrt{x^T A_{SS} x} = \begin{cases} A_{SS} x / \sqrt{x^T A_{SS} x}, & A_{SS} x \neq 0, \\ U_S \Lambda_S^{1/2} \{y : \|y\| \leq 1\}, & A_{SS} x = 0. \end{cases}$$

where  $U_S \in \mathbb{R}^{|S| \times |S|}$  has  $u_{i,S}$  as the  $i$ th columns and  $\Lambda_S^{1/2}$  is the diagonal matrix with  $\lambda_{i,S}^{1/2}, i = 1, \dots, |S|$ , as diagonal elements. For  $x$  with  $A_{SS} x \neq 0$ , there is a relationship between  $\partial \sqrt{x^T A_{SS} x}$  and  $\partial \sqrt{y^T A y}$  at  $y = g_S(x)$  given by

$$\begin{aligned} \left\{ \frac{Ay}{\sqrt{y^T A y}} \right\}_S &= \frac{A_{SS} x}{\sqrt{x^T A_{SS} x}}, \\ \left\{ \frac{Ay}{\sqrt{y^T A y}} \right\}_{S^c} &= \frac{A_{S^c S} x}{\sqrt{x^T A_{SS} x}}. \end{aligned}$$

Subgradients are bounded above:

$$\|\nabla h_{A,S}(x)\|_\infty \leq \|A_{SS}\|_2^{1/2} \leq \|A\|_2^{1/2}, \quad \forall \nabla f_{A,S}(x) \in \partial \sqrt{x^T A_{SS} x}.$$

*Proof.* (i) For  $x$  with  $Ax \neq 0$ ,  $f(x)$  is differentiable and the subgradient of  $f$  at  $x$  is simply the matrix derivative. By definition, for  $x$  with  $Ax = 0$ , the subgradient  $v$  of  $f$  at  $x$  satisfies the following inequality

$$\sqrt{y^T A y} \geq \langle y - x, v \rangle, \quad (22)$$

for all  $y$ . Choosing  $y = 2x$  and  $y = 0$  yield  $0 \geq \langle x, v \rangle$  and  $0 \geq -\langle x, v \rangle$ , implying  $\langle x, v \rangle = 0$ . The inequality (22) reduces to  $\sqrt{y^T A y} \geq \langle y, v \rangle$ , for any  $y$ . If  $Ay = 0$ , a similar argument implies that  $\langle y, v \rangle = 0$ . Hence  $v \perp y$  for every  $y$  with  $Ay = 0$ .

Let  $j_0$  be the smallest index such that  $\lambda_{j_0} > 0$ . Because  $u_j$ 's form an orthonormal basis, any arbitrary vector  $y$  can be written as  $y = \sum_{j=1}^p \beta_j u_j$ . Moreover, the null space of  $A$  is the span of  $u_1, \dots, u_{j_0-1}$ . Thus, the subgradient  $v$  can be written as  $v = \sum_{j=j_0}^p \alpha_j u_j$ . Thus, using the spectral decomposition of  $A$  as  $A = \sum_{j=j_0}^p \lambda_j u_j u_j^T$ , we can write  $f(y) = \{\sum_{j=j_0}^p \lambda_j \beta_j^2\}^{1/2}$ . On the other hand,  $\langle y, v \rangle = \sum_{j=j_0}^p \alpha_j \beta_j$ . Thus, the inequality (22) further reduces to

$$\left\{ \sum_{j=j_0}^p \lambda_j \beta_j^2 \right\}^{1/2} \geq \sum_{j=j_0}^p \alpha_j \beta_j, \quad \forall \beta_j \in \mathbb{R}.$$

It follows from the Cauchy-Schwartz inequality that the left hand side of the inequality is bounded from above;

$$\sum_{j=j_0}^p \alpha_j \beta_j = \sum_{j=j_0}^p \frac{\alpha_j}{\lambda_j^{1/2}} \lambda_j^{1/2} \beta_j \leq \left\{ \sum_{j=j_0}^p \frac{\alpha_j^2}{\lambda_j} \right\}^{1/2} \left\{ \sum_{j=j_0}^p \lambda_j \beta_j^2 \right\}^{1/2}.$$

Thus,

$$\partial f(x) = \left\{ v : v = \sum_{j=j_0}^p \alpha_j v_j, \sum_{j=j_0}^p \frac{\alpha_j^2}{\lambda_j} \leq 1, \alpha_j \in \mathbb{R} \right\}.$$

It is easy to see that this set is the image of the map  $U\Lambda^{1/2}$  on the closed ball of radius 1.

Given that  $\|x\|_\infty \leq \|x\|$ , to establish the bound in the  $\ell_\infty$ -norm, we compute the bound in the Euclidean norm. We use the same notation as in (i). For  $x$  with  $Ax \neq 0$ ,

$$\left\| \frac{Ax}{\sqrt{x^T A x}} \right\| = \frac{\|U\Lambda^{1/2}\Lambda^{1/2}U^T x\|}{\|\Lambda^{1/2}U^T x\|} \leq \|U\Lambda^{1/2}\|_2.$$

But  $\|U\Lambda^{1/2}\|_2 = \sup_{\|x\|=1} \|U\Lambda^{1/2}x\| = \sup_{x \in \mathbb{R}^K} \|U\Lambda^{1/2}(U^T x)\| / \|U^T x\| = \|A\|_2^{1/2}$ , because  $\|U^T x\| = \|x\|$ . For  $x$  with  $Ax = 0$ ,  $\|\Lambda^{1/2}y/\|y\|\| \leq \|A\|_2^{1/2}$  for

every  $y$ . Because of the form of the subdifferential and the fact that  $\|Ux\| = \|x\|$ , the result follows.

(ii) Let  $B_S$  be a product of elementary matrices for row and column exchange such that  $B_S g_S(x) = (x, 0)$ . Notice that  $B_S = B_S^{-1}$  and that  $B_S = B_S^T$  since  $B_S$  only rearranges elements of vectors and exchanges rows by multiplication from the left. Note also that  $\|B_S\|_2 \leq \|B_S\|_{\infty/\infty} = 1$ , since  $\|C\|_2 \leq \|C\|_{\infty/\infty}$  for  $C = C^T$  and each row of  $B_S$  has only one element with value 1. Because

$$\{h_{A,S}(x)\}^2 = g_S(x)^T A g_S(x) = (B_S g_S(x))^T (B_S A B_S) (B_S g_S(x)) = x^T A_{SS} x,$$

the subdifferential of  $h_{A,S}(x)$  follows from (ii). For  $x$  with  $A_{SS}x \neq x$  and  $y = g_S(x)$ ,  $Ay = B_S A B_S(x, 0)^T = B_S(A_{SS}x, A_{S^c S}^T x)^T \neq 0$  because of invertibility of  $B_S$ . The relationship holds since

$$\begin{bmatrix} (Ay/\sqrt{y^T Ay})_S \\ (Ay/\sqrt{y^T Ay})_{S^c} \end{bmatrix} = B_S \frac{Ay}{\sqrt{y^T Ay}} = \frac{1}{\sqrt{x^T A_{SS} x}} \begin{bmatrix} A_{SS} x \\ A_{S^c S} x \end{bmatrix}$$

An  $\ell_\infty$ -bound follows from (i) and the fact that  $\|A_{SS}\|_2 \leq \|B_S\|_2^2 \|A\|_2 = \|A\|_2$ .  $\square$

**Lemma 7.** For sample correlation matrices  $\hat{\Psi}_n = (\hat{\Psi}_n^{(1)}, \dots, \hat{\Psi}_n^{(K)})$  and any  $\rho_n > 0$ , the convex problem (3) has a unique solution  $\hat{\Theta}_{\rho_n} = (\hat{\Theta}_{\rho_n}^{(1)}, \dots, \hat{\Theta}_{\rho_n}^{(K)})$  with  $\hat{\Theta}_{\rho_n}^{(k)} > 0, k = 1, \dots, K$ , characterized by

$$n^{-1} n_k (\psi_{n,ij}^{(k)} - [\{\hat{\Theta}_{\rho_n}^{(k)}\}^{-1}]_{ij}) + \rho_n \hat{U}_{1,ij}^{(k)} + \rho_n \rho_2 \hat{U}_{2,ij}^{(k)} = 0, \quad (23)$$

with  $\hat{U}_{1,ij}^{(k)} \in \partial|\hat{\theta}_{\rho_n,ij}^{(k)}|$  and  $(\hat{U}_{2,ij}^{(1)}, \dots, \hat{U}_{2,ij}^{(K)})^T \in \partial\sqrt{\hat{\Theta}_{\rho_n,ij}^T L \hat{\Theta}_{\rho_n,ij}}$  for every  $i \neq j$  and  $k = 1, \dots, K$ . Moreover,

$$n^{-1} n_k (\psi_{n,ii}^{(k)} - [\{\hat{\Theta}_{\rho_n}^{(k)}\}^{-1}]_{ii}) + \rho_n \hat{U}_{1,ij}^{(k)} + \rho_n \rho_2 \hat{U}_{2,ij}^{(k)} = 0, \quad (24)$$

with  $\hat{U}_{1,ij}^{(k)} = \hat{U}_{2,ij}^{(k)} = 0$  for every  $i = 1, \dots, p$ , and  $k = 1, \dots, K$ .

For each  $(i, j) \in S$ , let  $S_{ij} = \{k : \Theta_{0,ij}^{(k)} \neq 0\}$ . The convex problem (21) has a unique solution  $\check{\Theta}_{\rho_n} = (\check{\Theta}_{\rho_n}^{(1)}, \dots, \check{\Theta}_{\rho_n}^{(K)})$  with  $\check{\Theta}_{\rho_n}^{(k)} > 0, k = 1, \dots, K$ , characterized by

$$n^{-1} n_k (\psi_{n,ij}^{(k)} - [\{\check{\Theta}_{\rho_n}^{(k)}\}^{-1}]_{ij}) + \rho_n \check{U}_{1,ij}^{(k)} + \rho_n \rho_2 \check{U}_{2,ij}^{(k)} = 0, \quad (25)$$

with  $\check{U}_{1,ij}^{(k)} \in \partial|\check{\theta}_{\rho_n,ij}^{(k)}|$  and  $\check{U}_{2,ij}^{(k)} \in \partial\sqrt{\{\check{\Theta}_{\rho_n,ij}\}_{S_{ij}}^T L_{S_{ij} S_{ij}} \{\check{\Theta}_{\rho_n,ij}\}_{S_{ij}}}$  for every  $i \neq j$  and  $k = 1, \dots, K$ . Moreover,

$$n^{-1} n_k (\psi_{n,ii}^{(k)} - [\{\check{\Theta}_{\rho_n}^{(k)}\}^{-1}]_{ii}) + \rho_n \check{U}_{1,ij}^{(k)} + \rho_n \rho_2 \check{U}_{2,ij}^{(k)} = 0, \quad (26)$$

with  $\check{U}_{1,ij}^{(k)} = \check{U}_{2,ij}^{(k)} = 0$  for every  $i = 1, \dots, p$ , and  $k = 1, \dots, K$ .

*Proof.* A proof for the uniqueness of the solution is similar to the proof of Lemma 3 of Ravikumar et al. [25]. The rest is the KKT condition using Lemma 6.  $\square$

We choose a pair  $\tilde{U} = (\tilde{U}_1, \tilde{U}_2)$  of the subgradients of the first and second regularization terms evaluated at  $\check{\Theta}_{\rho_n}$ . For each  $(i, j)$  with  $\Omega_{0,ij} = 0$  or with  $L\check{\Theta}_{\rho_n,ij} = 0$ , set

$$\tilde{U}_{1,ij}^{(k)} = \rho_n^{-1} n^{-1} n_k (-\psi_{n,ij}^{(k)} + [\{\check{\Theta}_{\rho_n}^{(k)}\}^{-1}]_{ij}), \quad \tilde{U}_{2,ij}^{(k)} = 0, \quad k = 1, \dots, K.$$

For  $(i, j)$  with  $\omega_{0,ij}^{(k)} \neq 0$ , for all  $k = 1, \dots, K$ , set

$$\tilde{U}_{1,ij}^{(k)} = \check{U}_{1,ij}^{(k)}, \quad \tilde{U}_{2,ij}^{(k)} = \check{U}_{2,ij}^{(k)}, \quad k = 1, \dots, K.$$

For  $(i, j)$  with  $L\check{\Theta}_{\rho_n,ij} \neq 0$ ,  $\Omega_{0,ij} \neq 0$  but  $\omega_{0,ij}^{(k')} = 0$  for some  $k'$ , set

$$\tilde{U}_{1,ij}^{(k)} = \rho_n^{-1} n^{-1} n_k (-\psi_{n,ij}^{(k)} + [\{\check{\Theta}_{\rho_n}^{(k)}\}^{-1}]_{ij}) - \rho_2 \frac{l_k \check{\Theta}_{\rho_n,ij}}{\sqrt{\check{\Theta}_{\rho_n,ij}^T L \check{\Theta}_{\rho_n,ij}}},$$

and

$$\tilde{U}_{2,ij}^{(k)} = \frac{l_k^T \check{\Theta}_{\rho_n,ij}}{\sqrt{\check{\Theta}_{\rho_n,ij}^T L \check{\Theta}_{\rho_n,ij}}},$$

if  $\omega_{0,ij}^{(k)} = 0$ . Otherwise, let

$$\tilde{U}_{1,ij}^{(k)} = \check{U}_{1,ij}^{(k)}, \quad \tilde{U}_{2,ij}^{(k)} = \frac{l_k^T \check{\Theta}_{\rho_n,ij}}{\sqrt{\check{\Theta}_{\rho_n,ij}^T L \check{\Theta}_{\rho_n,ij}}}.$$

Here,  $l_k$  is the  $k$ th row of  $L$ .

The main idea of the proof is to show that  $(\check{\Theta}_{\rho_n}, \tilde{U})$  satisfies the optimality conditions of the original problem with probability tending to 1. In particular, we show the following equation, which holds by construction of  $\tilde{U}_1$  and  $\tilde{U}_2$ , is in fact the KKT condition of the original problem (3):

$$n^{-1} n_k (\hat{\Psi}_n^{(k)} - \{\check{\Theta}_{\rho_n}^{(k)}\}^{-1}) + \rho_n \tilde{U}_1^{(k)} + \rho_n \rho_2 \tilde{U}_2^{(k)} = 0. \quad (27)$$

To this end, we show that  $\tilde{U}_1$  and  $\tilde{U}_2$  are both subgradients of the original problem. We can then conclude that the oracle estimator in the restricted problem (21) is the solution to the original problem (3). Then it follows from the uniqueness of the solution that  $\check{\Theta}_{\rho_n} = \hat{\Theta}_{\rho_n}$ .

Let  $\Xi^{(k)} = \hat{\Psi}_n^{(k)} - \Psi_0^{(k)}$ ,  $R^{(k)}(\Delta^{(k)}) = \{\check{\Theta}_{\rho_n}^{(k)}\}^{-1} - \Psi_0^{(k)} + \Psi_0^{(k)} \Delta^{(k)} \Psi_0^{(k)}$ , and  $\check{\Delta}^{(k)} = \check{\Theta}_{\rho_n}^{(k)} - \Theta_0^{(k)}$ .

**Lemma 8.** *Suppose that  $\max\{\|\Xi^{(k)}\|_\infty, \|R^{(k)}(\check{\Delta}^{(k)})\|_\infty\} \leq \alpha \rho_n / 8$ , and  $\rho_2 \leq \alpha^2 / \{4\|L\|_2^{1/2}(2 - \alpha)\}$ . Suppose moreover that  $L\check{\Theta}_{\rho_n,ij} \neq 0$  for  $(i, j) \in S$ . Then  $|\tilde{U}_{1,ij}^{(k)}| < 1$  for  $(i, j) \in (S^{(k)})^c$ .*

*Proof.* We rewrite (27) to obtain

$$\frac{n_k}{n} \Psi_0^{(k)} \check{\Delta}^{(k)} \Psi_0^{(k)} + \frac{n_k}{n} \Xi^{(k)} - \frac{n_k}{n} R^{(k)}(\check{\Delta}^{(k)}) + \rho_n \tilde{U}_1^{(k)} + \rho_n \rho_2 \tilde{U}_2^{(k)} = 0.$$

We further rewrite the above equation via vectorization;

$$\frac{n_k}{n} (\Psi_0^{(k)} \otimes \Psi_0^{(k)}) \vec{\Delta}^{(k)} + \frac{n_k}{n} \vec{\Xi}^{(k)} - \frac{n_k}{n} \vec{R}^{(k)}(\check{\Delta}^{(k)}) + \rho_n \vec{U}_1^{(k)} + \rho_n \rho_2 \vec{U}_2^{(k)} = 0.$$

We separate this equation into two equations depending on  $S^{(k)}$ ;

$$\begin{aligned} \frac{n_k}{n} \Gamma_{S^{(k)} S^{(k)}}^{(k)} \vec{\Delta}_{S^{(k)}}^{(k)} + \frac{n_k}{n} \vec{\Xi}_{S^{(k)}}^{(k)} - \frac{n_k}{n} \vec{R}_{S^{(k)}}^{(k)}(\check{\Delta}^{(k)}) + \rho_n \vec{U}_{1, S^{(k)}}^{(k)} + \rho_n \rho_2 \vec{U}_{2, S^{(k)}}^{(k)} &= 0, \\ \frac{n_k}{n} \Gamma_{(S^{(k)})^c S^{(k)}}^{(k)} \vec{\Delta}_{(S^{(k)})^c}^{(k)} + \frac{n_k}{n} \vec{\Xi}_{(S^{(k)})^c}^{(k)} - \frac{n_k}{n} \vec{R}_{(S^{(k)})^c}^{(k)}(\check{\Delta}^{(k)}) + \\ \rho_n \vec{U}_{1, (S^{(k)})^c}^{(k)} + \rho_n \rho_2 \vec{U}_{2, (S^{(k)})^c}^{(k)} &= 0. \end{aligned} \quad (28)$$

where  $(\vec{U}_l)_J \equiv \vec{U}_{k, J}, l = 1, 2$ . Here we used  $\check{\Delta}_{(S^{(k)})^c}^{(k)} = 0$ . Since  $\Gamma_{S^{(k)} S^{(k)}}^{(k)}$  is invertible, we solve the first equation to obtain

$$\frac{n_k}{n} \vec{\Delta}_{S^{(k)}}^{(k)} = (\Gamma_{S^{(k)} S^{(k)}}^{(k)})^{-1} \left\{ -\frac{n_k}{n} \vec{\Xi}_{S^{(k)}}^{(k)} + \frac{n_k}{n} \vec{R}_{S^{(k)}}^{(k)}(\check{\Delta}^{(k)}) - \rho_n \vec{U}_{1, S^{(k)}}^{(k)} - \rho_n \rho_2 \vec{U}_{2, S^{(k)}}^{(k)} \right\}.$$

Substituting this expression into (28) yields

$$\begin{aligned} \vec{U}_{1, (S^{(k)})^c}^{(k)} &= \rho_n^{-1} \Gamma_{(S^{(k)})^c S^{(k)}}^{(k)} (\Gamma_{S^{(k)} S^{(k)}}^{(k)})^{-1} \left( \frac{n_k}{n} \vec{\Xi}_{S^{(k)}}^{(k)} - \frac{n_k}{n} \vec{R}_{S^{(k)}}^{(k)}(\check{\Delta}^{(k)}) \right) \\ &\quad + \Gamma_{(S^{(k)})^c S^{(k)}}^{(k)} (\Gamma_{S^{(k)} S^{(k)}}^{(k)})^{-1} \vec{U}_{1, S^{(k)}}^{(k)} + \rho_2 \Gamma_{(S^{(k)})^c S^{(k)}}^{(k)} (\Gamma_{S^{(k)} S^{(k)}}^{(k)})^{-1} \vec{U}_{2, S^{(k)}}^{(k)} \\ &\quad - \rho_n^{-1} \left( \frac{n_k}{n} \vec{\Xi}_{(S^{(k)})^c}^{(k)} - \frac{n_k}{n} \vec{R}_{(S^{(k)})^c}^{(k)}(\check{\Delta}^{(k)}) \right) - \rho_2 \vec{U}_{2, (S^{(k)})^c}^{(k)}. \end{aligned}$$

Taking the  $\ell_\infty$ -norm yields

$$\begin{aligned} \left\| \vec{U}_{1, (S^{(k)})^c}^{(k)} \right\|_\infty &\leq \rho_n^{-1} \left\| \Gamma_{(S^{(k)})^c S^{(k)}}^{(k)} (\Gamma_{S^{(k)} S^{(k)}}^{(k)})^{-1} \right\|_{\infty/\infty} (\|\vec{\Xi}_{S^{(k)}}^{(k)}\|_\infty + \|\vec{R}_{S^{(k)}}^{(k)}(\check{\Delta}^{(k)})\|_\infty) \\ &\quad + \left\| \Gamma_{(S^{(k)})^c S^{(k)}}^{(k)} (\Gamma_{S^{(k)} S^{(k)}}^{(k)})^{-1} \right\|_{\infty/\infty} (\|\vec{U}_{1, S^{(k)}}^{(k)}\|_\infty + \rho_2 \|\vec{U}_{2, S^{(k)}}^{(k)}\|_\infty) \\ &\quad + \rho_n^{-1} (\|\vec{\Xi}_{(S^{(k)})^c}^{(k)}\|_\infty + \|\vec{R}_{(S^{(k)})^c}^{(k)}(\check{\Delta}^{(k)})\|_\infty) + \rho_2 \|\vec{U}_{2, (S^{(k)})^c}^{(k)}\|_\infty \\ &\leq \frac{2 - \alpha}{\rho_n} (\|\vec{\Xi}_{(S^{(k)})^c}^{(k)}\|_\infty + \|\vec{R}_{(S^{(k)})^c}^{(k)}(\check{\Delta}^{(k)})\|_\infty) \\ &\quad + 1 - \alpha + (2 - \alpha) \rho_2 \|L\|_2^{1/2}. \end{aligned}$$

Here we used that  $\|Ax\|_\infty \leq \|A\|_{\infty/\infty} \|x\|_\infty$  and  $\|\Gamma_{(S^{(k)})^c S^{(k)}}^{(k)} (\Gamma_{S^{(k)} S^{(k)}}^{(k)})^{-1}\|_{\infty/\infty} \leq 1 - \alpha$ , and applied Lemma 6 to bound  $\|\vec{U}_{2, (S^{(k)})^c}^{(k)}\|_\infty$  and  $\|\vec{U}_{2, S^{(k)}}^{(k)}\|_\infty$  by  $\|L\|_2^{1/2}$ . We also used  $\|\vec{U}_{1, S^{(k)}}^{(k)}\|_\infty = \|\vec{U}_{1, S^{(k)}}^{(k)}\|_\infty \leq 1$  by construction of  $\tilde{U}_1$  and the assumption that  $\check{\Theta}_{\rho_n}^{(k)} \neq 0$  for  $(i, j) \in S^{(k)}$ . It follows by the assumption of the

lemma that

$$\begin{aligned}\|\tilde{U}_{(S^{(k)})^c}^{(k)}\|_\infty &\leq \frac{2-\alpha}{\rho_n} \frac{\alpha\rho_n}{4} + (1-\alpha) + (2-\alpha)\rho_2\|L\|_2^{1/2} \\ &\leq 1 - \frac{\alpha}{2} - \frac{\alpha^2}{4} + \frac{\alpha^2}{4} < 1.\end{aligned}\quad \square$$

**Lemma 9** (Lemma 5 of Ravikumar et al. [25]). *Suppose that  $\|\Delta\|_\infty \leq 1/(3\kappa_\Psi d)$  with*

$$(\Delta^{(k)})_{(S^{(k)} \cup \{(i,i): i=1,\dots,p\})^c} = 0.$$

*Then  $\|H^{(k)}\|_{\infty/\infty} \leq 3/2$  where  $H^{(k)} \equiv \sum_{j=1}^\infty (-1)^j (\Psi_0^{(k)} \Delta^{(k)})^j$ ,  $k = 1, \dots, K$ , and  $R^{(k)}(\Delta^{(k)})$  has representation  $R^{(k)}(\Delta^{(k)}) = \Psi_0^{(k)} \Delta^{(k)} \Psi_0 \Delta H^{(k)} \Psi_0^{(k)}$  with  $\|R^{(k)}(\Delta^{(k)})\|_\infty \leq (3/2)d\|\Delta^{(k)}\|_\infty^2(\kappa_\Psi)^3$ .*

**Lemma 10.** *Suppose  $\|\Delta\|_2 \leq 1/(2\max_k \|\Psi_0^{(k)}\|_2)$  with  $\Delta_{(S^{(k)} \cup \{(i,i): i=1,\dots,p\})^c}^{(k)} = 0$ .*

*Then  $\|H^{(k)}\|_{\infty/\infty} \leq 2$  where  $H^{(k)} \equiv \sum_{t=1}^\infty (-1)^t (\Psi_0^{(k)} \Delta^{(k)})^t$ ,  $k = 1, \dots, K$ , and  $R^{(k)}(\Delta^{(k)})$  has representation  $R^{(k)}(\Delta^{(k)}) = \Psi_0^{(k)} \Delta^{(k)} \Psi_0 \Delta H^{(k)} \Psi_0^{(k)}$  with  $\|R^{(k)}(\Delta^{(k)})\|_\infty \leq 2\lambda_\Theta^3 \|\Delta^{(k)}\|_2^2$ .*

*Proof.* Note that the Neumann series for a matrix  $(I - A)^{-1}$  converges if the operator norm of  $A$  is strictly less than 1, and that the  $\ell_\infty$ -norm is bounded by the operator norm. A proof is similar to that of Lemma 5 of Ravikumar et al. [25] with the induced infinity norm  $\|\cdot\|_{\infty/\infty}$  replaced by the operator norm in appropriate inequalities.  $\square$

The following lemma is similar to the statement of Lemma 6 of Ravikumar et al. [25].

**Lemma 11.** *Suppose that*

$$r \equiv \frac{4}{\min_k \pi_k} \kappa_\Gamma (\max_k \|\Xi^{(k)}\|_\infty + \rho_n + \rho_n \rho_2 \|L\|_2^{1/2}) < \frac{1}{6d \max\{\kappa_\Psi, \kappa_\Psi^3 \kappa_\Gamma\}},$$

*for  $k = 1, \dots, K$ . Suppose moreover that  $(\Theta_0^{(k)} \otimes \Theta_0^{(k)})_{S^{(k)} S^{(k)}}$  are invertible for  $k = 1, \dots, K$ . Then with probability  $1 - 2K \exp(-n \min_k \pi_k^2/2)$ ,*

$$\max_k \|\tilde{\Theta}_{\rho_n}^{(k)} - \Theta_0^{(k)}\|_\infty \leq (3/2)r.$$

*Proof.* We apply Shauder's fixed point theorem on the event  $\min_k \pi_k/2 \leq n_k/n$ , which holds with probability  $1 - 2K \exp(-n \min_k \pi_k^2/2)$  by Lemma 5 with  $\epsilon = \min_k \pi_k/2$ . We first define the function  $f_k$  and its domain  $\mathcal{D}_k$  to which the fixed point theorem applies. Let  $\bar{S}^{(k)} = S^{(k)} \cup \{(i,i) : 1 \leq i \leq p\}$ , and define

$$\mathcal{D}_k = \{A \in \mathbb{S}^{p \times p} : x^T(A + \Theta_0^{(k)})x \geq 0, \forall x \in \mathbb{R}^p, \|A_{\bar{S}^{(k)}}\|_\infty \leq r, A_{(\bar{S}^{(k)})^c} = 0\},$$

where  $\mathbb{S}^{p \times p}$  is the space of symmetric  $p \times p$  matrices. Then  $\mathcal{D}_k$  is a convex, compact subset of the set of  $\mathbb{S}^{p \times p}$ .



Let  $\check{U}_l^{(k)} \in \mathbb{R}^{p \times p}$ ,  $l = 1, 2$ , be zero-filled matrices whose  $(i, j)$ -element is  $\check{U}_{l,ij}^{(k)}$  in Lemma 7 if  $(i, j) \in S^{(k)}$  and zero otherwise. Define the map  $g_k$  on the set of invertible matrices in  $\mathbb{R}^{p \times p}$  by  $g_k(B) = (n_k/n)(B^{-1} - \hat{\Psi}_n^{(k)}) - \rho_n \check{U}_1^{(k)} - \rho_n \rho_2 \check{U}_2^{(k)}$ . Note that  $\{g_k(\check{\Theta}_{\rho_n}^{(k)})\}_{S^{(k)}} = 0$  is the KKT condition for the restricted problem (21). Let  $\delta > 0$  be a constant such that  $\delta < \min\{1/2, 1/\{10(4dr + 1)\}\}r$  and  $\delta + r \leq 1/\{6d \max\{\kappa_\Psi, \kappa_\Psi^3 \kappa_\Gamma\}\}$ . Define a continuous function  $f_k : \mathcal{D}_k \mapsto \mathcal{D}_k$  as

$$(f_k(A))_{ij} = \begin{cases} A_{ij} & i = j, \\ \left\{ h_k(A) \Theta_0^{(k)} g_k(A + \Theta_0^{(k)} + \delta I) \Theta_0^{(k)} + A \right\}_{ij} & i \neq j, (i, j) \in S^{(k)} \\ 0, & \text{otherwise,} \end{cases}$$

where

$$h_k(A) \equiv \frac{2^{-1} \min\{\lambda_1(A + \Theta_0^{(k)}), 2^{-1}\} + 2^{-1}}{\max\{|\lambda_1(\{\Theta_0^{(k)} g_k(A + \Theta_0^{(k)} + \delta I) \Theta_0^{(k)}\}_{S^{(k)}} - I)|, 1\}}.$$

Let  $\tilde{f}_k(A) = h_k(A) \Theta_0^{(k)} g_k(A + \Theta_0^{(k)} + \delta I) \Theta_0^{(k)}$ . Then  $f_k(A) = (\tilde{f}_k(A))_{S^{(k)}} + A$  for  $A \in \mathcal{D}_k$ .

We now verify the conditions of Shauder's fixed point theorem below. Once these conditions are established, the theorem yields that  $f_k(A) = A$ . Since  $(f_k(A))_{(\bar{S}^{(k)})^c} = A$  for any  $A \in \mathcal{D}_k$ , and  $h_k(A) > 0$ , the solution  $A$  to  $f_k(A) = A$  is determined by  $(\Theta_0^{(k)} g_k(A + \Theta_0^{(k)} + \delta I) \Theta_0^{(k)})_{S^{(k)}} = 0$ . Vectorizing this equation to obtain  $(\Theta_0^{(k)} \otimes \Theta_0^{(k)})_{S^{(k)} S^{(k)}} \{g_k(A + \Theta_0^{(k)} + \delta I)\}_{S^{(k)}} = 0$ , it follows from the invertibility of  $(\Theta_0^{(k)} \otimes \Theta_0^{(k)})_{S^{(k)} S^{(k)}}$  that  $\{g_k(A + \Theta_0^{(k)} + \delta I)\}_{S^{(k)}} = 0$ . By the uniqueness of the KKT condition, the solution is  $A = \check{\Theta}_{\rho_n}^{(k)} - \Theta_0^{(k)} - \delta I$ . Since  $A \in \mathcal{D}_k$ , and  $\delta < r/2$ , we conclude  $\|\check{\Theta}_{\rho_n}^{(k)} - \Theta_0^{(k)}\|_\infty \leq (3/2)r$ .

In the following, we write  $\vec{A} = \text{vec}(A)$  for a matrix  $A$  for notational convenience. For  $J \subset \{(i, j) : i, j = 1, \dots, p\}$ ,  $\text{vec}(A)_J$  should be understood as  $\vec{A}_J$ .

The function  $f_k$  is continuous on  $\mathcal{D}_k$ . To see this, note first that  $A + \Theta_0^{(k)} + \delta I$  is positive definite for every  $A \in \mathcal{D}_k$  so that the inversion is continuous. Note also that all elements in the matrices involved with eigenvalues in  $h_k(A)$  are uniformly bounded in  $\mathcal{D}_k$ , and hence the eigenvalues are also uniformly bounded.

To show that  $f_k(A) \in \mathcal{D}_k$ , first we show that  $f_k(A) + \Theta_0^{(k)}$  is positive semidefinite. This follows because for any  $x \in \mathbb{R}^p$

$$\begin{aligned} & x^T (f_k(A) + \Theta_0^{(k)}) x \\ &= x^T \{(\tilde{f}_k(A))_{S^{(k)}} - I\} x + x^T (A + \Theta_0^{(k)}) x + x^T x \\ &\geq h_k(A) \lambda_1(\{\Theta_0^{(k)} g_k(A + \Theta_0^{(k)} + \delta I) \Theta_0^{(k)}\}_{S^{(k)}} - I) \|x\|^2 \\ &\quad + \lambda_1(A + \Theta_0^{(k)}) \|x\|^2 + \|x\|^2 \geq 0. \end{aligned}$$

To see this, note that if  $\lambda_A \equiv \lambda_1(\{\Theta_0^{(k)} g_k(A + \Theta_0^{(k)} + \delta I) \Theta_0^{(k)}\}_{S^{(k)}} - I)$  is positive,

then the inequality easily follows. On the other hand, if  $\lambda_A < -1$ , we have

$$\begin{aligned} h_k(A)\lambda_A\|x\|^2 &\geq -2^{-1}\min\{\lambda_1(A + \Theta_0^{(k)}), 2^{-1}\}\|x\|^2 - 2^{-1}\|x\|^2 \\ &\geq -(\lambda_1(A + \Theta_0^{(k)})/2 + 1/2)\|x\|^2. \end{aligned}$$

Lastly, if  $-1 \leq \lambda_A < 0$ , we have

$$\begin{aligned} h_k(A)\lambda_A\|x\|^2 &\geq -|\lambda_A|[2^{-1}\min\{\lambda_1(A + \Theta_0^{(k)}), 2^{-1}\} + 1/2]\|x\|^2 \\ &\geq -|\lambda_A|(\lambda_1(A + \Theta_0^{(k)})/2 + 1/2)\|x\|^2. \end{aligned}$$

Next, we show that  $\|f_k(A)_{\bar{S}^{(k)}}\|_\infty \leq r$ . Because  $\text{diag}(f_k(A)) = \text{diag}(A)$ , it suffices to show  $\|f_k(A)_{S^{(k)}}\|_\infty \leq r$ . Since  $\delta + r \leq 1/\{6d \max\{\kappa_\Psi, \kappa_\Psi^3 \kappa_\Gamma\}\}$ ,

$$\|\Psi_0^{(k)}(A + \delta I)\|_{\infty/\infty} \leq \kappa_\Psi d \|A + \delta I\|_\infty \leq \kappa_\Psi d(r + \delta) \leq 1/3.$$

It then follows from Lemma 9 that

$$\begin{aligned} R(A + \delta I) &= \left(A + \delta I + \Theta_0^{(k)}\right)^{-1} - \Psi_0^{(k)} + \Psi_0^{(k)}(A + \delta I)\Psi_0^{(k)} \\ &= \{\Psi_0^{(k)}(A + \delta I)\}^2 H^{(k)} \Psi_0^{(k)}. \end{aligned}$$

Thus, adding and subtracting  $\Psi_0^{(k)}$  yields

$$\begin{aligned} \tilde{f}_k(A) + A &= h_k(A)\Theta_0^{(k)}((n_k/n)\{\Psi_0^{(k)}(A + \delta I)\}^2 H^{(k)} \Psi_0^{(k)} - (n_k/n)\Xi^{(k)} - \rho_n \check{U}_1^{(k)} \\ &\quad - \rho_n \rho_2 \check{U}_2^{(k)})\Theta_0^{(k)} + (1 - (n_k/n)h_k(A))A - (n_k/n)\delta h_k(A)I. \end{aligned}$$

Vectorization and restriction on  $S^{(k)}$  gives

$$\begin{aligned} \|\text{vec}(f_k(A))_{S^{(k)}}\|_\infty &= \|\text{vec}(\tilde{f}_k(A) + A)_{S^{(k)}}\|_\infty \\ &\leq (n_k/n)h_k(A)\|\{(\Gamma^{(k)})^{-1}\}_{S^{(k)}S^{(k)}}\text{vec}(\{\Psi_0^{(k)}(A + \delta I)\}^2 H^{(k)} \Psi_0^{(k)})_{S^{(k)}}\|_\infty \\ &\quad + (1 - (n_k/n)h_k(A))\|\text{vec}(A)_{S^{(k)}}\|_\infty + (n_k/n)\delta \\ &\quad + h_k(A)\|\{(\Gamma^{(k)})^{-1}\}_{S^{(k)}S^{(k)}}\{\text{vec}((n_k/n)\Xi^{(k)})_{S^{(k)}} \\ &\quad + \rho_n \text{vec}(\check{U}_1^{(k)})_{S^{(k)}} + \rho_n \rho_2 \text{vec}(\check{U}_2^{(k)})_{S^{(k)}}\}\|_\infty, \end{aligned} \tag{29}$$

where  $\{(\Gamma^{(k)})^{-1}\}_{S^{(k)}S^{(k)}} = (\Theta^{(k)} \otimes \Theta_0^{(k)})_{S^{(k)}S^{(k)}}$ . Here we used  $h_k(A) \leq (1/4 + 1/2)/1 = 3/4$ . For the first term of the upper bound in (29), it follows by the inequality  $\|Ax\|_\infty \leq \|A\|_{\infty/\infty}\|x\|_\infty$  for  $A \in \mathbb{R}^{p \times p}$  and  $x \in \mathbb{R}^p$ , Lemma 9 and the choice of  $\delta$  satisfying  $\delta + r \leq 1/\{6d \max\{\kappa_\Psi, \kappa_\Psi^3 \kappa_\Gamma\}\}$  that

$$\begin{aligned} &\|\{(\Gamma^{(k)})^{-1}\}_{S^{(k)}S^{(k)}}\text{vec}(\{\Psi_0^{(k)}(A + \delta I)\}^2 H^{(k)} \Psi_0^{(k)})_{S^{(k)}}\|_\infty \\ &\leq \kappa_\Gamma \|R^{(k)}(A + \delta I)\|_\infty \leq \kappa_\Gamma \frac{3}{2}d \|A + \delta I\|_\infty^2 \kappa_\Psi^3 \leq \kappa_\Gamma \frac{3}{2}d \|A + \delta I\|_\infty (r + \delta) \kappa_\Psi^3 \\ &\leq (r + \delta)/4. \end{aligned}$$

For the second term, it follows by the assumption, the inequality that  $\|Ax\|_\infty \leq \|A\|_{\infty/\infty} \|x\|_\infty$  for  $A \in \mathbb{R}^{p \times p}$  and  $x \in \mathbb{R}^p$ , and Lemma 6 that

$$\begin{aligned} & \|\{(\Gamma^{(k)})^{-1}\}_{S^{(k)}S^{(k)}} \left\{ \frac{n_k}{n} \text{vec}(\Xi^{(k)})_{S^{(k)}} + \rho_n \text{vec}(\check{U}_1^{(k)})_{S^{(k)}} + \rho_n \rho_2 \text{vec}(\check{U}_2^{(k)})_{S^{(k)}} \right\}\|_\infty \\ & \leq \kappa_\Gamma (\|\Xi^{(k)} + \rho_n + \rho_n \rho_2 \|L\|_2^{1/2}) = (\min_k \pi_k) r/4 \leq (n_k/n) r/2. \end{aligned}$$

Thus, we can further bound  $\|\text{vec}((\tilde{f}_k(A) + A)_{S^{(k)}})\|_\infty$  by

$$\begin{aligned} & \frac{n_k}{n} h_k(A) \frac{r+\delta}{4} + \frac{n_k}{n} h_k(A) \frac{r}{2} + \left(1 - \frac{n_k}{n} h_k(A)\right) r + \frac{n_k}{n} \delta = \\ & r \left\{1 - \frac{n_k}{n} \frac{h_k(A)}{4}\right\} + \frac{n_k}{n} \left\{1 + \frac{h_k(A)}{4}\right\} \delta. \end{aligned} \quad (30)$$

Since

$$\begin{aligned} & \|(\Theta_0^{(k)} g_k(A + \Theta_0^{(k)} + \delta I) \Theta_0^{(k)})_{S^{(k)}}\|_\infty \\ & \leq \|A_{S^{(k)}}\|_\infty + \|\Theta_0^{(k)} g_k(A + \Theta_0^{(k)} + \delta I) \Theta_0^{(k)} + A\|_{S^{(k)}}\|_\infty, \end{aligned}$$

and  $\delta \leq r/2$ , a similar reasoning shows that

$$\begin{aligned} & \|(\Theta_0^{(k)} g_k(A + \Theta_0^{(k)} + \delta I) \Theta_0^{(k)})_{S^{(k)}}\|_\infty \\ & \leq (n_k/n) \|\{(\Gamma^{(k)})^{-1}\}_{S^{(k)}S^{(k)}} \text{vec}\left(\left\{\Psi_0^{(k)}(A + \delta I)\right\}^2 H^{(k)} \Psi_0^{(k)}\right)_{S^{(k)}}\|_\infty \\ & \quad + (2 - (n_k/n)) \|\text{vec}(A)_{S^{(k)}}\|_\infty + (n_k/n) \delta \\ & \quad + \|\{(\Gamma^{(k)})^{-1}\}_{S^{(k)}S^{(k)}} \left\{ (n_k/n) \text{vec}(\Xi^{(k)})_{S^{(k)}} + \rho_n \text{vec}(\check{U}_1^{(k)})_{S^{(k)}} \right. \\ & \quad \left. + \rho_n \rho_2 \text{vec}(\check{U}_2^{(k)})_{S^{(k)}} \right\}\|_\infty \\ & \leq \frac{r+\delta}{4} + \frac{r}{2} + 2r + \delta \leq 4r. \end{aligned}$$

Thus, the inequality  $\|B\|_2 \leq \|B\|_{\infty/\infty}$  for  $B = B^T$  implies that

$$\begin{aligned} & |\lambda_1(\{\Theta_0^{(k)} g_k(A + \Theta_0^{(k)} + \delta I) \Theta_0^{(k)}\}_{S^{(k)}} - I)| \\ & \leq \|\lambda_1(\{\Theta_0^{(k)} g_k(A + \Theta_0^{(k)} + \delta I) \Theta_0^{(k)}\}_{S^{(k)}})\|_2 + 1 \\ & \leq \|\lambda_1(\{\Theta_0^{(k)} g_k(A + \Theta_0^{(k)} + \delta I) \Theta_0^{(k)}\}_{S^{(k)}})\|_{\infty/\infty} + 1 \\ & \leq 4dr + 1. \end{aligned}$$

Hence  $h_k(A) \geq 1/(8dr + 2)$  for every  $A \in \mathcal{D}_k$ .

Now (30) is further bounded by  $r$ :

$$r \left\{1 - \frac{n_k}{n} \frac{h_k(A)}{4}\right\} + \frac{n_k}{n} \left\{1 + \frac{h_k(A)}{4}\right\} \delta$$

$$\begin{aligned}
&\leq r \left\{ 1 - \frac{n_k}{n} \frac{h_k(A)}{4} \right\} + \frac{n_k}{n} \left\{ 1 + \frac{h_k(A)}{4} \right\} \frac{r}{10(4dr+1)} \\
&\leq r \left\{ 1 - \frac{n_k}{n} \frac{h_k(A)}{4} \right\} + \frac{n_k}{n} \left\{ 1 + \frac{h_k(A)}{4} \right\} \frac{h_k(A)r}{5} \\
&\leq r - \frac{n_k}{n} \frac{h_k(A) - h_k^2(A)}{20} r \leq r.
\end{aligned}$$

Here we used the fact that  $\delta \leq r/\{10(4dr+1)\}$  and  $1/(8dr+2) \leq h_k(A) < 1$ . Thus,  $\|(f_k(A))_{S^{(k)}}\|_\infty \leq r$ .

Since  $(f_k(A))_{(S^{(k)})^c} = 0$  by definition, all the conditions for the fixed point theorem are established. This completes the proof.  $\square$

We are now ready to prove Theorem 2. Note that Condition 7 implies that

$$\rho_n < \min \left\{ \frac{\min_k \pi_k}{72d\kappa_\Gamma} \min \left\{ \frac{1}{\kappa_\Psi}, \frac{1}{\kappa_\Psi^3 \kappa_\Gamma}, \frac{\min_k \pi_k}{56\kappa_\Psi^3 \kappa_\Gamma} \alpha \right\}, \frac{c_8}{6}, \frac{c_9 \min_k \sqrt{d_k}}{12} \right\}.$$

*Proof of Theorem 2.* We prove that the oracle estimator  $\check{\Theta}_{\rho_n}$  satisfies (I) the model selection consistency and (II) the KKT conditions of the original problem (3) with  $(\check{\Theta}_{\rho_n}, \tilde{U}_1, \tilde{U}_2)$ . The model selection consistency of  $\hat{\Theta}_{\rho_n} = \check{\Theta}_{\rho_n}$  then follows by the uniqueness of the solution to the original problem. The following discussion is on the event that  $\min_k \pi_k/2 \leq n_k/n, k = 1, \dots, K$ , and  $\max_k \|\Xi^{(k)}\|_\infty \leq \alpha/8$ . Note that this event has probability approaching 1 by Lemmas 4 and 5.

First we obtain an  $\ell_\infty$ -bound of the error of the oracle estimator. Note that by Condition 7 and the fact that  $\alpha \in [0, 1]$

$$\frac{\alpha}{8} + 1 + \rho_2 \|L\|_2^{1/2} \leq \frac{\alpha}{8} + 1 + \frac{\alpha^2}{4(2-\alpha)} \leq 3.$$

Thus, it follows from Condition 7 that

$$\begin{aligned}
\frac{4}{\min_k \pi_k} \kappa_\Gamma \left( \|\Xi^{(k)}\|_\infty + \rho_n + \rho_n \rho_2 \|L\|_2^{1/2} \right) &< \frac{12\kappa_\Gamma}{\min_k \pi_k} \frac{\min_k \pi_k}{72d\kappa_\Gamma} \min \left\{ \frac{1}{\kappa_\Psi}, \frac{1}{\kappa_\Psi^3 \kappa_\Gamma} \right\} \\
&= \frac{1}{6d \max\{\kappa_\Psi, \kappa_\Psi^3 \kappa_\Gamma\}}.
\end{aligned}$$

Because  $(\Theta_0^{(k)} \otimes \Theta_0^{(k)})_{S^{(k)} S^{(k)}}$  is invertible by Condition 5, we can apply Lemma 11 to obtain  $\|\check{\Theta}_{\rho_n}^{(k)} - \Theta_0^{(k)}\|_\infty \leq (6/\min_k \pi_k) \kappa_\Gamma (\|\Xi^{(k)}\|_\infty + \rho_n + \rho_n \rho_2 \|L\|_2^{1/2})$  with probability approaching 1.

As a consequence of the  $\ell_\infty$ -bound,  $\check{\Theta}_{\rho_n, ij} \neq 0$  for  $(i, j) \in S$ , because  $\|\check{\Theta}_{\rho_n}^{(k)} - \Theta_0^{(k)}\|_\infty \leq 3\rho_n \leq c_8/2 < \min_{k=1, \dots, K, i \neq j} |\theta_{0, ij}^{(k)}|$  by Conditions 6 and 7. This establishes the model selection consistency of the oracle estimator.

Next, we show that the Oracle estimator satisfies the KKT condition of the original problem (3). As the first step, we prove  $\tilde{U}_{1, ij}^{(k)} \in \partial \check{\Theta}_{\rho_n}^{(k)}$  for every  $i, j, k$  with probability approaching 1. Since  $\check{\Theta}_{\rho_n, ij} \neq 0$  for  $(i, j) \in S$  with probability

approaching 1,  $\tilde{U}_{1,ij}^{(k)} = \check{U}_{1,ij}^{(k)}$  for  $(i, j) \in S^{(k)}$  by construction. For  $(i, j) \in (S^{(k)})^c$ , we need to prove  $|\tilde{U}_{1,ij}^{(k)}| < 1$  for every  $i, j, k$ . To this end, it suffices to verify that  $\|R^{(k)}(\check{\Theta}_{\rho_n}^{(k)} - \Theta_0^{(k)})\|_\infty \leq \alpha/8$  and apply Lemma 8. Applying Lemma 9 with  $\|\check{\Theta}_{\rho_n}^{(k)} - \Theta_0^{(k)}\|_\infty \leq (6/\min_k \pi_k) \kappa_\Gamma (\|\Xi^{(k)}\|_\infty + \rho_n + \rho_n \rho_2 \|L\|_2^{1/2})$  and Condition 7 gives

$$\begin{aligned} \|R^{(k)}(\check{\Theta}_{\rho_n}^{(k)} - \Theta_0^{(k)})\|_\infty &\leq \frac{3}{2} d\kappa_\Psi^3 \|\check{\Theta}_{\rho_n}^{(k)} - \Theta_0^{(k)}\|_\infty^2 \leq \frac{3}{2} d\kappa_\Psi^3 \frac{324\kappa_\Gamma^2}{\min_k \pi_k^2} \rho_n^2 \\ &\leq \frac{486d\kappa_\Psi^3 \kappa_\Gamma^2}{\min_k \pi_k^2} \left\{ \frac{\min_k \pi_k}{72d\kappa_\Gamma} \frac{\min_k \pi_k}{56\kappa_\Psi^3 \kappa_\Gamma} \alpha \right\} \rho_n \leq \frac{\alpha}{8}. \end{aligned}$$

Next, we prove that  $\tilde{U}_{2,ij} \in \partial\sqrt{\check{\Theta}_{\rho_n,ij} L \check{\Theta}_{\rho_n,ij}}$  for every  $(i, j)$ . For  $(i, j)$  with  $\omega_{0,ij}^{(k)} \neq 0$  for all  $k = 1, \dots, K$ ,  $\tilde{U}_{2,ij} = \check{U}_{\rho_n} \in \partial\sqrt{\check{\Theta}_{\rho_n,ij} L \check{\Theta}_{\rho_n,ij}}$ . For  $(i, j)$  with  $\Omega_{0,ij} = 0$ ,  $\tilde{U}_{2,ij} = 0 \in \partial\sqrt{\check{\Theta}_{\rho_n,ij} L \check{\Theta}_{\rho_n,ij}}$  by Lemma 6. For  $(i, j)$  with  $\Omega_{0,ij} \neq 0$  and  $\omega_{0,ij}^{(k')} = 0$  for some  $k'$ ,

$$\tilde{U}_{2,ij} = L \check{\Theta}_{\rho_n,ij} / \sqrt{\check{\Theta}_{\rho_n,ij} L \check{\Theta}_{\rho_n,ij}} \in \partial\sqrt{\check{\Theta}_{\rho_n,ij} L \check{\Theta}_{\rho_n,ij}}$$

if  $L \check{\Theta}_{\rho_n,ij} \neq 0$ . To see  $L \check{\Theta}_{\rho_n,ij} \neq 0$  holds with probability approaching 1, let  $(k, k') \in S$  with  $k \neq k'$  such that  $\Theta_{0,ij}^{(k)}/\sqrt{d_k} - \Theta_{0,ij}^{(k')}/\sqrt{d_{k'}} \neq 0$ . This pair  $(k, k')$  exists by Condition 6 and the assumption  $L\Theta_{0,ij} \neq 0$ . We assume without loss of generality  $\theta_{0,ij}^{(k)}/\sqrt{d_k} - \theta_{0,ij}^{(k')}/\sqrt{d_{k'}} > 0$ . Since  $\|\check{\Theta}_{\rho_n}^{(k)} - \Theta_0^{(k)}\|_\infty \leq 3\rho_n \leq c_9 \min_k \sqrt{d_k}/12$ , it follows from Condition 7 that

$$\begin{aligned} \frac{\check{\theta}_{\rho_n,ij}^{(k)}}{\sqrt{d_k}} - \frac{\check{\theta}_{\rho_n,ij}^{(k')}}{\sqrt{d_{k'}}} &\geq \frac{\theta_{0,ij}^{(k)}}{\sqrt{d_k}} - \frac{\theta_{0,ij}^{(k')}}{\sqrt{d_{k'}}} - 3\rho_n \left( \frac{1}{\sqrt{d_k}} + \frac{1}{\sqrt{d_{k'}}} \right) \\ &\geq c_9 - 3\rho_n \left( \max_{W_{k,k'} \neq 0} \frac{1}{\sqrt{d_k}} + \frac{1}{\sqrt{d_{k'}}} \right) \geq \frac{1}{2} c_9. \end{aligned}$$

Hence,  $\check{\Theta}_{\rho_n,ij}^T L \check{\Theta}_{\rho_n,ij} \geq W_{kk'} c_9^2/4 > 0$  or  $L \check{\Theta}_{\rho_n,ij} \neq 0$ .

Finally, we show that Equation (27) for the KKT condition holds. For the  $(i, j)$ -element of the equation with  $\Omega_{0,ij} = 0$ , this equation holds by construction for every  $k = 1, \dots, K$ . For the  $(i, j)$ -element with  $\omega_{0,ij}^{(k)} \neq 0$  for every  $k = 1, \dots, K$ , the equation holds for every  $k = 1, \dots, K$ , because it is the equation for the KKT condition of the corresponding element in a restricted problem (21). For  $(i, j)$ -element with  $\Omega_{0,ij} \neq 0$  and  $\omega_{0,ij}^{(k')} = 0$  for some  $k'$ , note that  $\check{\Theta}_{\rho_n,ij} \neq 0$  with probability approaching 1 and that the rearrangement in  $\Theta_{ij}$  and corresponding exchange of rows and columns of  $L$  for each  $i, j$  does not change the original and restricted optimization problems (3) and (21). Thus, with the appropriate rearrangement of elements and exchange of rows and columns,  $\tilde{U}_{2,ij}^{(k)}$  with  $\omega_{0,ij}^{(k)} \neq 0$  is in fact  $\check{U}_{2,ij}^{(k)}$ . Thus for such  $k$  the equation

holds because of the corresponding KKT condition in the restricted problem (21). For other  $k$ , the equation holds by construction. We thus conclude the oracle estimator satisfies the KKT condition of the original problem (3). This completes the proof.  $\square$

*Proof of Corollary 1.* In the proof of Theorem 2, the  $\ell_\infty$ -bound of the error yields

$$\|\hat{\Theta}_{\rho_n}^{(k)} - \Theta_0^{(k)}\|_\infty = O_P(\kappa_\Gamma \rho_n).$$

Note that if one of two matrices  $A$  and  $B$  is diagonal,  $\|AB\|_\infty \leq \|A\|_\infty \|B\|_\infty$ . Thus, we can proceed in the same way as in the proof of Theorem 2 of Rothman et al. [26] to conclude that

$$\|\hat{\Omega}_n^{(k)} - \Omega_0^{(k)}\|_\infty = O_P(\kappa_\Gamma \rho_n).$$

The result follows from a similar argument to the proof of Corollary 3 in Ravikumar et al. [25].  $\square$

*Proof of Corollary 2.* It follows from Condition 8 and Lemma 1 applied to  $\check{\Theta}_{\rho_n}$  that  $\|\check{\Theta}_{\rho_n}^{(k)} - \Theta_0^{(k)}\|_2 \leq 1/(2\lambda_\Theta)$ . Then we can apply Lemma 10 instead of Lemma 9. The rest is similar to the proof of Theorem 2.  $\square$

### Hierarchical clustering

For simplicity, we prove Theorem 3 for the case of  $K = 2$ ; the proof can be easily generalized to  $K > 2$ . Let  $X$  and  $Y$  be the random variable from the first and subpopulation, respectively. Suppose that  $X = (X_1, \dots, X_p)^T \sim N(\mu_X, \Sigma_X)$  with  $\mu_X = (\mu_{1,X}, \dots, \mu_{p,X})$  and the spectral decomposition  $\Sigma_X = Q_X \Lambda_X Q_X^T$  of  $\Sigma_X$  where  $\lambda_{1,X}, \dots, \lambda_{p,X}$  are the eigenvalues of  $\Sigma_X$  and that  $Y \sim N(\mu_Y, \Sigma_Y)$  with  $\mu_Y = (\mu_{1,Y}, \dots, \mu_{p,Y})$  and the spectral decomposition  $\Sigma_Y = Q_Y \Lambda_Y Q_Y^T$  of  $\Sigma_Y$  where  $\lambda_{1,Y}, \dots, \lambda_{p,Y}$  are the eigenvalues of  $\Sigma_Y$ . Define  $Z = (X - Y) = (Z_1, \dots, Z_p)^T \sim N(\mu_Z, \Sigma_Z)$  with  $\mu_Z = (\mu_{1,Z}, \dots, \mu_{p,Z})$  and the spectral decomposition  $\Sigma_Z = Q_Z \Lambda_Z Q_Z^T$  of  $\Sigma_Z$  where  $\lambda_{1,Z}, \dots, \lambda_{p,Z}$  are the eigenvalues of  $\Sigma_Z$ . Let  $\tilde{X} = (\tilde{X}_1, \dots, \tilde{X}_p)^T = \Lambda_X^{1/2} Q_X^T \Sigma_X^{-1/2} X$ ,  $\tilde{Y} = (\tilde{Y}_1, \dots, \tilde{Y}_p)^T = \Lambda_Y^{1/2} Q_Y^T \Sigma_Y^{-1/2} Y$  and  $\tilde{Z} = (\tilde{Z}_1, \dots, \tilde{Z}_p)^T = \Lambda_Z^{1/2} Q_Z^T \Sigma_Z^{-1/2} Z$ . Then  $\tilde{X} \sim N(\tilde{\mu}_X, \Lambda_X)$ ,  $\tilde{Y} \sim N(\tilde{\mu}_Y, \Lambda_Y)$  and  $\tilde{Z} \sim N(\tilde{\mu}_Z, \Lambda_Z)$ , where

$$\begin{aligned} \tilde{\mu}_X &= (\tilde{\mu}_{1,X}, \dots, \tilde{\mu}_{p,X})^T \equiv \Lambda_X^{1/2} Q_X^T \Sigma_X^{-1/2} \mu_X, \\ \tilde{\mu}_Y &= (\tilde{\mu}_{1,Y}, \dots, \tilde{\mu}_{p,Y})^T \equiv \Lambda_Y^{1/2} Q_Y^T \Sigma_Y^{-1/2} \mu_Y, \\ \tilde{\mu}_Z &= (\tilde{\mu}_{1,Z}, \dots, \tilde{\mu}_{p,Z})^T \equiv \Lambda_Z^{1/2} Q_Z^T \Sigma_Z^{-1/2} \mu_Z. \end{aligned}$$

Let also

$$\begin{aligned} \mu_X^2 &= \|\tilde{\mu}_X^2\|/p, \quad \mu_Y^2 = \|\tilde{\mu}_Y^2\|/p, \quad \mu_Z^2 = \|\tilde{\mu}_Z^2\|/p, \\ \bar{\lambda}_X &= \sum_{k=1}^p \lambda_{k,X}/p, \quad \bar{\lambda}_Y = \sum_{k=1}^p \lambda_{k,Y}/p, \quad \bar{\lambda}_Z = \sum_{k=1}^p \lambda_{k,Z}/p. \end{aligned}$$

**Lemma 12** (Lemma 1 of Borysov et al. [1]). *Let  $W_1, \dots, W_p$  be independent non-negative random variables with finite second moments. Let  $S = \sum_{j=1}^p (W_j - \mathbb{E} W_j)$  and  $v = \sum_{j=1}^p \mathbb{E} W_j^2$ . Then for any  $t > 0$   $P(S \leq -t) \leq \exp(-t^2/(2v))$ .*

The following lemma is an extension of Lemma 2 in Borysov et al. [1].

**Lemma 13.** *Let  $0 < a < \mu_{\tilde{X}}^2 + \bar{\lambda}_X$ . Then*

$$P(\|X\|^2 < ap) \leq \exp \left( -\frac{p^2(\mu_{\tilde{X}}^2 + \bar{\lambda}_X - a)^2}{2 \sum_{j=1}^p (\tilde{\mu}_{j,X}^4 + 6\tilde{\mu}_{j,X}^2 \lambda_{j,X} + 3\lambda_{j,X}^2)} \right).$$

*Proof.* Note that elements of  $\tilde{X}$  are independent and that  $\tilde{X}_j \sim N(\tilde{\mu}_{j,X}, \lambda_{j,X})$ . Thus, we have

$$\begin{aligned} \mathbb{E} \tilde{X}_j^2 &= \tilde{\mu}_{j,X}^2 + \lambda_{j,X}, \quad \text{Var}(\tilde{X}_j^2) = 2(\lambda_{j,X}^2 + 2\tilde{\mu}_{j,X}^2 \lambda_{j,X}), \\ \mathbb{E} \tilde{X}_j^4 &= \tilde{\mu}_{j,X}^4 + 6\tilde{\mu}_{j,X}^2 \lambda_{j,X} + 3\lambda_{j,X}^2. \end{aligned}$$

Applying Lemma 12 with  $W_i = \tilde{X}_i^2$ , since  $P(\|X\|^2 < ap) = P(\|\tilde{X}\|^2 < ap)$ , we get

$$\begin{aligned} P(\|X\|^2 < ap) &= P \left[ \sum_{j=1}^p (\tilde{X}_j^2 - \tilde{\mu}_{j,X}^2 - \lambda_{j,X}) < -p(\mu_{\tilde{X}}^2 + \bar{\lambda}_X - a) \right] \\ &\leq \exp \left( -\frac{p^2(\mu_{\tilde{X}}^2 + \bar{\lambda}_X - a)^2}{2 \sum_{j=1}^p (\tilde{\mu}_{j,X}^4 + 6\tilde{\mu}_{j,X}^2 \lambda_{j,X} + 3\lambda_{j,X}^2)} \right). \quad \square \end{aligned}$$

The following is an extension of Lemma 3 in Borysov et al. [1].

**Lemma 14.** *Let  $a > \bar{\lambda}_X + \mu_{\tilde{X}}^2$ . Then*

$$P(\|X\|^2 > ap) \leq \exp \left( -\frac{1}{2} \left( p + \sum_{j=1}^p \frac{a}{\lambda_{j,X}} - \sum_{j=1}^p \sqrt{1 + 2\frac{a}{\lambda_{j,X}}} \right) \right).$$

*Proof.* By Markov's inequality, for  $t > \sum_{j=1}^p \lambda_{j,X} + \tilde{\mu}_{j,X}^2$ , we get

$$\begin{aligned} P \left( \sum_{j=1}^p X_j^2 \geq t \right) &= P \left( \sum_{j=1}^p \tilde{X}_j^2 \geq t \right) \\ &= P \left[ \exp \left( \sum_{j=1}^p \gamma \tilde{X}_j^2 - \gamma \lambda_{j,X} - \gamma \tilde{\mu}_{j,X}^2 \right) \geq \exp \left( \gamma t - \gamma \sum_{j=1}^p (\lambda_{j,x} + \tilde{\mu}_{j,X}^2) \right) \right] \\ &\leq \exp \left( -\gamma \left( t - \sum_{j=1}^p (\tilde{\mu}_{j,X}^2 + \lambda_{j,X}) \right) \right) \prod_{j=1}^p \mathbb{E} \exp((\gamma \lambda_{j,X}) \tilde{X}_j^2 / \lambda_{j,x}) \end{aligned}$$

$$= \exp \left( -\gamma \left( t - \sum_{j=1}^p \tilde{\mu}_{j,X}^2 \right) \right) \prod_{j=1}^p \exp \left( -\gamma \lambda_{j,X} - \frac{1}{2} \log(1 - 2\gamma \lambda_{j,X}) \right) \times \exp \left( \frac{\gamma \tilde{\mu}_{j,X}^2}{1 - 2\gamma \lambda_{j,X}} \right).$$

Since for all  $u \in (0, 1)$ ,  $-\log(1-u)-u \leq u^2/\{2(1-u)\}$  (see page 28 of Boucheron et al. [2]), the above display is bounded above by

$$\exp \left( -\gamma \left( t - \sum_{i=1}^p \tilde{\mu}_{i,X}^2 \right) \right) \prod_{i=1}^p \exp \left( \frac{\gamma^2 \lambda_{i,X}^2}{1 - 2\gamma \lambda_{i,X}} \right) \exp \left( \frac{\gamma \tilde{\mu}_{i,X}^2}{1 - 2\gamma \lambda_{i,X}} \right).$$

Using the following result from Boucheron et al. [2]

$$\inf_{\gamma \in (0, 1/c)} \frac{v\gamma^2}{2(1-c\gamma)} - t\gamma = -\frac{v}{c^2} h \left( \frac{ct}{v} \right).$$

wherein  $h(u) = 1 + u - \sqrt{1 + 2u}$ ,  $u > 0$ , we further obtain the upper bound

$$\exp \left( \gamma \sum_{i=1}^p \tilde{\mu}_{i,X}^2 \right) \prod_{i=1}^p \exp \left( -\frac{1}{2} \left( 1 + \frac{t}{\lambda_{i,X}p} - \sqrt{1 + \frac{2t}{\lambda_{i,X}p}} \right) \right) \exp \left( \frac{\gamma \tilde{\mu}_{i,X}^2}{1 - 2\gamma \lambda_{i,X}} \right).$$

Taking  $\gamma \downarrow 0$ , the upper bound becomes

$$\exp \left( -\frac{1}{2} \left( p + \sum_{i=1}^p \frac{t}{\lambda_{i,X}p} - \sum_{i=1}^p \sqrt{1 + \frac{2t}{\lambda_{i,X}p}} \right) \right).$$

Choosing  $t = ap$ , we have

$$P \left( \sum_{i=1}^p \tilde{X}_i^2 \geq ap \right) \leq \exp \left( -\frac{1}{2} \left( p + \sum_{i=1}^p \frac{a}{\lambda_{i,X}} - \sum_{i=1}^p \sqrt{1 + 2\frac{a}{\lambda_{i,X}}} \right) \right).$$

Note that  $f(u) = (1 + 2u)^{1/2} \leq u$  for  $u \geq 0$  because  $f'(0) = 1$  and  $f'$  is decreasing for  $u > 0$ . Thus,  $P \left( \sum_{i=1}^p \tilde{X}_i^2 \geq ap \right) \rightarrow 0$  as  $p \rightarrow \infty$ .  $\square$

*Proof of Theorem 3.* For simplicity, we present the proof for the case of  $K = 2$ ; the proof can be easily generalized to  $K > 2$ . Let  $n_1$  and  $n_2$  be the sample sizes for the first and second subpopulations, respectively. Define

$$\begin{aligned} E_1 &= \left\{ \max_{i,j} \|X_i - X_j\| < \min_{k,l} \|X_k - Y_l\| \right\}, & E_3 &= \left\{ \max_{i,j} \|X_i - X_j\|^2 < ap \right\}, \\ E_2 &= \left\{ \max_{i,j} \|Y_i - Y_j\| < \min_{k,l} \|X_k - Y_l\| \right\}, & E_4 &= \left\{ \max_{i,j} \|Y_i - Y_j\|^2 < ap \right\}, \\ E_5 &= \left\{ \max_{k,l} \|X_k - Y_l\|^2 > ap \right\}, \end{aligned}$$

for a fixed  $a > 0$  satisfying the assumption. The intersection  $E_1 \cap E_2$  is contained in the event that the clustering performs in the way that two subpopulations



are joined in the last step. The intersection  $E_3 \cap E_4 \cap E_5$  is also contained in  $E_1 \cap E_2$ , or in other words,  $P((E_1 \cap E_2)^c) \leq P(E_3^c) + P(E_4^c) + P(E_5^c)$ . Thus, it suffices to show that  $P(E_3^c) + P(E_4^c) + P(E_5^c) \rightarrow 0$  as  $n, p \rightarrow \infty$ .

For  $E_3^c$  and  $E_4^c$  we have by Lemma 14 that

$$\begin{aligned} P(E_3^c) &\leq \sum_{i,j}^n P(\|X_i - X_j\|^2 > ap) = \frac{n_1(n_1 - 1)}{2} P(\|X_1 - X_2\|^2 > ap) \\ &\leq \frac{n_1(n_1 - 1)}{2} \exp \left( -\frac{1}{2} \left( p + \sum_{l=1}^p \frac{a}{2\lambda_{l,X}} - \sum_{l=1}^p \sqrt{1 + \frac{a}{\lambda_{l,X}}} \right) \right) \\ &\leq \exp \left( -\frac{1}{2} \left( p + \sum_{l=1}^p \frac{a}{2\lambda_{l,X}} - \sum_{l=1}^p \sqrt{1 + \frac{a}{\lambda_{l,X}}} \right) + 2 \log n_1 \right) \\ &= \exp \left( -\frac{p}{2} \left( 1 + \frac{1}{p} \sum_{l=1}^p \frac{a}{2\lambda_{l,X}} - \frac{1}{p} \sum_{l=1}^p \sqrt{1 + \frac{a}{\lambda_{l,X}}} + 4 \frac{\log n_1}{p} \right) \right) \end{aligned}$$

and that

$$P(E_4^c) \leq \exp \left( -\frac{p}{2} \left( 1 + \frac{1}{p} \sum_{l=1}^p \frac{a}{2\lambda_{l,Y}} - \frac{1}{p} \sum_{l=1}^p \sqrt{1 + \frac{a}{\lambda_{l,Y}}} + 4 \frac{\log n_2}{p} \right) \right).$$

for  $a$  satisfying  $a > 2 \max\{\bar{\lambda}_X, \bar{\lambda}_Y\}$ .

Note that  $\log n_k/p \rightarrow 0, k = 1, 2$  as  $n_1, n_2, p \rightarrow \infty$ . Moreover  $x - \sqrt{1 + 2x} \geq 0$  for  $x > 0$ . Thus,  $P(E_3^c) \rightarrow 0$  and  $P(E_4^c) \rightarrow 0$  as  $n_1, n_2, p \rightarrow \infty$ . For  $E_5^c$ , we have by Lemma 13 that

$$\begin{aligned} P(E_5^c) &\leq \sum_{i,j} P(\|X_i - Y_j\|^2 < ap) \leq n_1 n_2 P(\|X_1 - Y_1\|^2 < ap) \\ &\leq \exp \left( -\frac{p^2(\mu_Z^2 + \bar{\lambda}_Z - a)^2}{2 \sum_{l=1}^p (\tilde{\mu}_{l,Z}^4 + 6\tilde{\mu}_{l,Z}^2 \lambda_{l,Z} + 3\lambda_{l,Z}^2)} + \log n_1 n_2 \right) \end{aligned}$$

for  $a < \mu_Z^2 + \bar{\lambda}_Z$ . Given the assumption  $c_{10} \leq \lambda_{j,X} \leq c_{11}, c_{10} \leq \lambda_{j,Y} \leq c_{11}, \max\{|\mu_{j,X}|, |\mu_{j,Y}|\} \leq c_{11}, j = 1, 2, \dots$ . Thus, we get  $P(E_5^c) \rightarrow 0$  as  $n_1, n_2, p \rightarrow \infty$ .

Since  $2\bar{\lambda}_X - \lambda_{p,X} - \lambda_{p,Y} \geq 2\bar{\lambda}_X - \bar{\lambda}_Z$ , and  $2\bar{\lambda}_Y - \lambda_{p,X} - \lambda_{p,Y} \geq 2\bar{\lambda}_Y - \bar{\lambda}_Z$ , the assumption that  $\mu_Z^2 > 2 \min\{\bar{\lambda}_X, \bar{\lambda}_Y\} - \lambda_{p,X} - \lambda_{p,Y}$  implies that there exists  $a$  such that  $a < \bar{\mu}_Z + \bar{\lambda}_Z$  and  $a > 2 \max\{\bar{\lambda}_X, \bar{\lambda}_Y\}$ . This completes the proof.  $\square$

## References

- [1] Petro Borysov, Jan Hannig, and JS Marron. Asymptotics of hierarchical clustering for growing dimension. *Journal of Multivariate Analysis*, 124:465–479, 2014. [MR3147338](#)
- [2] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford University Press, 2013. [MR3185193](#)

- [3] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, 3(1):1–122, 2011.
- [4] Tony Cai, Weidong Liu, and Xi Luo. A constrained  $\ell_1$  minimization approach to sparse precision matrix estimation. *J. Amer. Statist. Assoc.*, 106(494):594–607, 2011. ISSN 0162-1459. [MR2847973](#)
- [5] Fan RK Chung. *Spectral graph theory*, volume 92. American Mathematical Soc., 1997. [MR1421568](#)
- [6] Patrick Danaher, Pei Wang, and Daniela M Witten. The joint graphical lasso for inverse covariance estimation across multiple classes. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(2):373–397, 2014. [MR3164871](#)
- [7] Alexandre d’Aspremont, Onureena Banerjee, and Laurent El Ghaoui. First-order methods for sparse covariance selection. *SIAM J. Matrix Anal. Appl.*, 30(1):56–66, 2008. ISSN 0895-4798. [MR2399568](#)
- [8] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441, 2007.
- [9] Jian Guo, Elizaveta Levina, George Michailidis, and Ji Zhu. Joint estimation of multiple graphical models. *Biometrika*, 98(1):1–15, 2011. ISSN 0006-3444. [MR2804206](#)
- [10] Jian Huang, Shuangge Ma, Hongzhe Li, and Cun-Hui Zhang. The sparse Laplacian shrinkage estimator for high-dimensional regression. *Ann. Statist.*, 39(4):2021–2046, 2011. ISSN 0090-5364. [MR2893860](#)
- [11] Trey Ideker and Nevan J Krogan. Differential network biology. *Molecular systems biology*, 8(1), 2012.
- [12] Göran Jönsson, Johan Staaf, Johan Vallon-Christersson, Markus Ringnér, Karolina Holm, Cecilia Hegardt, Haukur Gunnarsson, Rainer Fagerholm, Carina Strand, Bjarni A Agnarsson, et al. Genomic subtypes of breast cancer identified by array-comparative genomic hybridization display distinct molecular and clinical characteristics. *Breast Cancer Research*, 12(3):1–14, 2010.
- [13] Mladen Kolar, Le Song, and Eric P Xing. Sparsistent learning of varying-coefficient models with structural changes. In *Advances in Neural Information Processing Systems*, pages 1006–1014, 2009.
- [14] Steffen L Lauritzen. *Graphical models*. Oxford University Press, 1996. [MR1419991](#)
- [15] Caiyan Li and Hongzhe Li. Variable selection and regression analysis for graph-structured covariates with an application to genomics. *Ann. Appl. Stat.*, 4(3):1498–1516, 2010. ISSN 1932-6157. [MR2758338](#)
- [16] Fan Li and Nancy R Zhang. Bayesian variable selection in structured high-dimensional covariate spaces with applications in genomics. *Journal of the American Statistical Association*, 105(491):1202–1214, 2010. [MR2752615](#)
- [17] F Liu, AC Lozano, S Chakraborty, and F Li. A graph laplacian prior for variable selection and grouping. *Biometrika*, 98(1):1–31, 2011.

- [18] Fei Liu, Sounak Chakraborty, Fan Li, Yan Liu, Aurelie C Lozano, et al. Bayesian regularization via graph laplacian. *Bayesian Analysis*, 9(2):449–474, 2014. [MR3217003](#)
- [19] Nicolai Meinshausen and Peter Bühlmann. High-dimensional graphs and variable selection with the lasso. *Ann. Statist.*, 34(3):1436–1462, 2006. ISSN 0090-5364. [MR2278363](#)
- [20] Sahand N. Negahban, Pradeep Ravikumar, Martin J. Wainwright, and Bin Yu. A unified framework for high-dimensional analysis of  $m$ -estimators with decomposable regularizers. *Stat. Sci.*, 27(4):538–557, 2012. [MR3025133](#)
- [21] Sahand N. Negahban, Pradeep Ravikumar, Martin J. Wainwright, and Bin Yu. Supplementary material for “a unified framework for high-dimensional analysis of  $m$ -estimators with decomposable regularizers”. *Stat. Sci.*, 2012. [MR3025133](#)
- [22] Charles M Perou, Therese Sørlie, Michael B Eisen, Matt van de Rijn, Stefanie S Jeffrey, Christian A Rees, Jonathan R Pollack, Douglas T Ross, Hilde Johnsen, Lars A Akslen, et al. Molecular portraits of human breast tumours. *Nature*, 406(6797):747–752, 2000.
- [23] Christine Peterson, Francesco C Stingo, and Marina Vannucci. Bayesian inference of multiple gaussian graphical models. *Journal of the American Statistical Association*, 110(509):159–174, 2015. [MR3338494](#)
- [24] Franck Rapaport, Andrei Zinovyev, Marie Dutreix, Emmanuel Barillot, and Jean-Philippe Vert. Classification of microarray data using gene networks. *BMC Bioinformatics*, 8, 2007.
- [25] Pradeep Ravikumar, Martin J. Wainwright, Garvesh Raskutti, and Bin Yu. High-dimensional covariance estimation by minimizing  $\ell_1$ -penalized log-determinant divergence. *Electron. J. Stat.*, 5:935–980, 2011. ISSN 1935-7524. [MR2836766](#)
- [26] Adam J. Rothman, Peter J. Bickel, Elizaveta Levina, and Ji Zhu. Sparse permutation invariant covariance estimation. *Electron. J. Stat.*, 2:494–515, 2008. ISSN 1935-7524. [MR2417391](#)
- [27] Nafiseh Sedaghat, Takumi Saegusa, Timothy Randolph, and Ali Shojaie. Comparative study of computational methods for reconstructing genetic networks of cancer-related pathways. *Cancer Informatics*, 13(Suppl 2):55–66, 09 2014.
- [28] Ali Shojaie and George Michailidis. Penalized principal component regression on graphs for analysis of subnetworks. In John D. Lafferty, Christopher K. I. Williams, John Shawe-Taylor, Richard S. Zemel, and Aron Culotta, editors, *NIPS*, pages 2155–2163. Curran Associates, Inc., 2010.
- [29] Nicolas Städler, Peter Bühlmann, and Sara Van De Geer.  $\ell_1$ -penalization for mixture regression models. *Test*, 19(2):209–256, 2010. [MR2677722](#)
- [30] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):267–288, 1996. [MR1379242](#)
- [31] Yu-Xiang Wang, James Sharpnack, Alex Smola, and Ryan J Tibshirani. Trend filtering on graphs. *arXiv preprint arXiv:1410.7690*, 2014.

- [32] Kilian Q Weinberger, Fei Sha, Qihui Zhu, and Lawrence K Saul. Graph laplacian regularization for large-scale semidefinite programming. In *Advances in neural information processing systems (NIPS)*, pages 1489–1496, 2006.
- [33] Ming Yuan. High dimensional inverse covariance matrix estimation via linear programming. *J. Mach. Learn. Res.*, 11:2261–2286, 2010. ISSN 1532-4435. [MR2719856](#)
- [34] Ming Yuan and Yi Lin. Model selection and estimation in the Gaussian graphical model. *Biometrika*, 94(1):19–35, 2007. ISSN 0006-3444. [MR2367824](#)
- [35] Peng Zhao and Bin Yu. On model selection consistency of lasso. *The Journal of Machine Learning Research*, 7:2541–2563, 2006. [MR2274449](#)
- [36] Peng Zhao, Guilherme Rocha, and Bin Yu. The composite absolute penalties family for grouped and hierarchical variable selection’. *Annals of Statistics*, 37(6A):3468–3497, 2009. [MR2549566](#)
- [37] Sen Zhao and Ali Shojaie. A significance test for graph-constrained estimation. *Biometrics* (forthcoming), 2015.