Learning High-Dimensional Differential Graphs From Multi-Attribute Data

Jitendra K. Tugnait

Abstract—We consider the problem of estimating differences in two Gaussian graphical models (GGMs) which are known to have similar structure. The GGM structure is encoded in its precision (inverse covariance) matrix. In many applications one is interested in estimating the difference in two precision matrices to characterize underlying changes in conditional dependencies of two sets of data. Existing methods for differential graph estimation are based on single-attribute (SA) models where one associates a scalar random variable with each node. In multiattribute (MA) graphical models, each node represents a random vector. In this paper, we analyze a group lasso penalized Dtrace loss function approach for differential graph learning from multi-attribute data. An alternating direction method of multipliers (ADMM) algorithm is presented to optimize the objective function. Theoretical analysis establishing consistency in support recovery and estimation in high-dimensional settings is provided. Numerical results based on synthetic as well as real data are presented.

Index Terms—Sparse graph learning; differential graph estimation; undirected graph; multi-attribute graphs.

I. INTRODUCTION

■ RAPHICAL models provide a powerful tool for analyzing multivariate data [1], [2]. In a statistical graphical model, the conditional statistical dependency structure among p random variables x_1, x_1, \dots, x_p , is represented using an undirected graph $\mathcal{G} = (V, \mathcal{E})$, where $V = \{1, 2, \dots, p\} = [p]$ is the set of p nodes corresponding to the p random variables x_i s, and $\mathcal{E} \subseteq V \times V$ is the set of undirected edges describing conditional dependencies among the components of x. The graph \mathcal{G} then is a conditional independence graph (CIG) where there is no edge between nodes i and j (i.e., $\{i, j\} \notin \mathcal{E}$) iff x_i and x_j are conditionally independent given the remaining p-2 variables x_{ℓ} , $\ell \in [p]$, $\ell \neq i$, $\ell \neq j$. In particular, Gaussian graphical models (GGMs) are CIGs where x is multivariate Gaussian. Suppose x has positive-definite covariance matrix Σ with inverse covariance matrix $\Omega = \Sigma^{-1}$. Then Ω_{ij} , the (i, j)-th element of Ω , is zero iff x_i and x_j are conditionally independent. Such models for x have been extensively studied. Given n samples of x, in high-dimensional settings where $p \gg 1$ and/or n is of the order of p, one estimates Ω under some sparsity constraints; see [3]-[6].

More recently there has been increasing interest in differential network analysis where one is interested in estimating the difference in two inverse covariance matrices [7]–[9]. Given observations x and y from two groups of subjects,

This work was supported by the National Science Foundation Grants ECCS-2040536 and CCF-2308473.

one is interested in the difference $\Delta = \Omega_y - \Omega_x$, where $\Omega_x = (E\{xx^{\top}\})^{-1}$ and $\Omega_y = (E\{yy^{\top}\})^{-1}$. The associated differential graph is $\mathcal{G}_{\Delta} = (V, \mathcal{E}_{\Delta})$ where $\{i,j\} \in \mathcal{E}_{\Delta}$ iff $\Delta_{ij} \neq 0$. It characterizes differences between the GGMs of the two sets of data. We use the term differential graph as in [10], [11] ([7]-[9]) use the term differential network). As noted in [9], in biostatistics, the differential network/graph describes the changes in conditional dependencies between components under different environmental or genetic conditions. For instance, one may be interested in the differences in the graphical models of healthy and impaired subjects, or models under different disease states, given gene expression data or functional MRI signals [3], [12], [13].

In the preceding graphs, each node represents a scalar random variable. In many applications, there may be more than one random variable associated with a node. This class of graphical models has been called multi-attribute (MA) graphical models in [14]-[17] and vector graphs or networks in [18]-[21]. In a gene regulatory network, one may have different molecular profiles available for a single gene, such as protein, DNA and RNA. Since these molecular profiles are on the same set of biological samples, they constitute multiattribute data for gene regulatory graphical models in [14], [16]. Consider p jointly Gaussian vectors $z_i \in \mathbb{R}^m$, $i \in [p]$. We associate z_i with the *i*th node of graph $\mathcal{G} = (V, \mathcal{E})$, $V = [p], \mathcal{E} \subseteq V \times V$. We now have m attributes per node. Now $\{i,j\} \in \mathcal{E}$ iff vectors z_i and z_j are conditionally independent given the remaining p-2 vectors $\{\boldsymbol{z}_{\ell}, \ell \in V \setminus \{i,j\}\}$. Let $\boldsymbol{x} = [\boldsymbol{z}_{1}^{\top} \ \boldsymbol{z}_{2}^{\top} \ \cdots \ \boldsymbol{z}_{p}^{\top}]^{\top} \in \mathbb{R}^{mp}$. Let $\boldsymbol{\Omega} = (E\{\boldsymbol{x}\boldsymbol{x}^{\top}\})^{-1}$ assuming $E\{xx^{\top}\} \succ 0$. Define the $m \times m$ subblock $\Omega^{(ij)}$ of Ω as $[\Omega^{(ij)}]_{rs} = [\Omega]_{(i-1)m+r,(j-1)m+s}, r,s=1,2,\cdots,m.$ Then we have the following equivalence [16, Sec. 2.1]

$$\{i, j\} \notin \mathcal{E} \Leftrightarrow \Omega^{(ij)} = \mathbf{0}.$$
 (1)

This paper is concerned with estimation of differential graphs from multi-attribute data. Given independent and identically distributed (i.i.d.) samples $\boldsymbol{x}(t), \ t=1,2,\cdots,n_x$, of $\boldsymbol{x}=[\boldsymbol{z}_1^\top \ \boldsymbol{z}_2^\top \ \cdots \ \boldsymbol{z}_p^\top]^\top \in \mathbb{R}^{mp}$ where $\boldsymbol{z}_i \in \mathbb{R}^m$, $i \in [p]$, are jointly Gaussian, and similarly given samples $\boldsymbol{y}(t)$, $t=1,2,\cdots,n_y$, of $\boldsymbol{y} \in \mathbb{R}^{mp}$, our objective is to estimate the difference $\boldsymbol{\Delta}=\boldsymbol{\Omega}_y-\boldsymbol{\Omega}_x$, and determine the differential graph $\mathcal{G}_{\Delta}=(V,\mathcal{E}_{\Delta})$ with edgeset $\mathcal{E}_{\Delta}=\{\{k,\ell\}: \|\boldsymbol{\Delta}^{(k\ell)}\|_F\neq 0\}$.

A. Related Work

All prior work on high-dimensional differential graph estimation from i.i.d. samples addresses single-attribute (SA) models where each node represents a scalar random variable. One naive approach would be to estimate the two precision

J.K. Tugnait is with the Department of Electrical & Computer Engineering, 200 Broun Hall, Auburn University, Auburn, AL 36849, USA. Email: tugnajk@auburn.edu.

matrices separately by any existing estimator (see [4], [5] and references therein) and then calculate their difference to estimate the differential graph. (This approach is also applicable to MA graphs.) This approach estimates twice the number of parameters, hence needs larger sample sizes for same accuracy, and also imposes sparsity constraints on each precision matrix for the methods to work. The same comment applies to methods such as [3], [6], where the two precision matrices and their differences are jointly estimated. A recent survey is in [22]. In these approaches, given $K \geq 2$ related groups of data, each p-variate and sharing the same set of nodes V, but possibly differing in connected edgesets, the objective is to jointly estimate the K precision matrices and their pairwise differences, with sparsity constraints on each of the K precision matrices and their pairwise differences. These approaches require each of the K precision matrices to be sparse. If only the differences in the precision matrices is of interest, alternative approaches exist where no sparsity constraints are imposed on individual precision matrices. For instance, direct estimation of the difference in the two precision matrices has been considered for SA graphs in [7]-[9], [12], [23]–[28], where only the difference is required to be sparse, not the two individual precision matrices. In [7]-[9], [23], [24], [28] precision difference matrix estimators are based on a D-trace loss [29], while [12] discusses a Dantzig selector type estimator. In [25]-[27] differential graph is estimated by directly modeling the ratio of the probability densities of the random vectors under the two graphs.

Estimation of MA differential graphs has not been investigated before. The work of [10], [11] is similar to an MA formulation except that in [10], [11], x(t) and y(t)are non-stationary ("functional" modeling), and instead of a single record (sample) of x(t), $t = 1, 2, \dots, n_x$ and y(t), $t = 1, 2, \dots, n_y$, as in this paper, they assume multiple independent observations of x(t), $t \in \mathcal{T}$ (a closed subset of real line), and y(t), $t \in \mathcal{T}$. The objective function in [11, Eqns. (10)-(11)] is the same as our objective function (3)-(4), but consequent estimation of edges and theoretical analysis are vastly different. We estimate edges as in (6), i.e., our threshold is set at zero and this is the method analyzed in our Theorem 1(iv) for graph recovery with high probability. In [10], [11], this threshold is set at a parameter $\epsilon_n > 0$ (see [11, Eqn. (13)]) which is a function of sample size n, number of nonzero entries in true Δ , smallest eigenvalues of true covariances Ω_{η}^{-1} and Ω_x^{-1} (in our notation), and several other factors. That is, ϵ_n is unknowable for practical implementation and it is used as a theoretical construct to establish graph support recovery in [11, Theorem 10]. In simulations, [10], [11] set $\epsilon_n = 0$. That is, [10], [11] do not analyze what they implement (the proof does not hold for $\epsilon_n = 0$), and they do not implement what they analyze (ϵ_n is unknowable). There is no counterpart to our Theorem 1 in [10], [11], and the methodology of our Theorem 1 allows us to set the edge detection threshold to zero. Our Theorem 2 follows the general framework of [30] to bound the Frobenius norm of the error in estimating Δ , and [10], [11] also follow the general framework of [30] for the same purpose. But their extension of this result to graph recovery does not permit zero threshold for edge detection. We attempt no such extension.

B. Our Contributions

In this paper, we analyze a group lasso penalized D-trace loss function approach for differential graph learning from MA data, extending the SA approach of [8], [28]. A two-block ADMM algorithm is presented to optimize the objective function. The two-block ADMM is guaranteed to be convergent unlike the three-block ADMM method used in [8]. Two different approaches to theoretical analysis of the proposed approach in high-dimensional settings are presented. Theorem 1 follows the approach(es) of [8], [16], [28], [29], [31] while Theorem 2 follows the general framework of [30], not used in [8], [28]. The general method of [31] requires an irrepresentability condition (see (20)) which is also required in [8], [28] for SA graphs, but is not needed by the method of [30], hence in our Theorem 2. Numerical results based on synthetic as well as real data are presented.

Preliminary version of parts of this paper appear in a conference paper [32]. Theorem 2, proof of Theorem 1 and real data example do not appear in [32].

C. Outline and Notation

The rest of the paper is organized as follows. A group lasso penalized D-trace loss function is presented in Sec. II for estimation of multi-attribute differential graph. An ADMM algorithm is presented in Sec. III to optimize the convex objective function. In Sec. IV we analyze the properties of the estimator of the difference $\Delta = \Omega_y - \Omega_x$. Theorem 1 follows the approach(es) of [8], [16], [28], [29], [31] while Theorem 2 follows the general framework of [30]. The general method of [31] requires an irrepresentability condition (see (20)) which is not needed by the method of [30]. On the other hand, our Theorem 2 does not have a result like Theorem 1(ii), the oracle property, nor does it have a result as in Theorem 1(iv), support recovery. Numerical results based on synthetic as well as real data are presented in Sec. V to illustrate the proposed approach. Proofs of Theorems 1 and 2 are given in Appendices A and B, respectively.

For a set V, |V| or $\operatorname{card}(V)$ denotes its cardinality. Given $A \in \mathbb{R}^{p \times p}$, we use $\phi_{\min}(A)$, $\phi_{\max}(A)$, |A| and $\operatorname{tr}(A)$ to denote the minimum eigenvalue, maximum eigenvalue, determinant and trace of A, respectively. For $B \in \mathbb{R}^{p \times q}$, we define $\|B\| = \sqrt{\phi_{\max}(B^{\top}B)}$, $\|B\|_F = \sqrt{\operatorname{tr}(B^{\top}B)}$, $\|B\|_1 = \sum_{i,j} |B_{ij}|$, where B_{ij} is the (i,j)-th element of B (also denoted by $[B]_{ij}$), $\|B\|_{\infty} = \max_{i,j} |B_{ij}|$ and $\|B\|_{1,\infty} = \max_i \sum_j |B_{ij}|$. The symbols \otimes and \boxtimes denote Kronecker product and Tracy-Singh product [33], respectively. In particular, given block partitioned matrices $A = [A_{ij}]$ and $B = [B_{k\ell}]$ with submatrices A_{ij} and $B_{k\ell}$, Tracy-Singh product yields another block partitioned matrix $A \boxtimes B = [A_{ij} \boxtimes B]_{ij} = [[A_{ij} \otimes B_{k\ell}]_{k\ell}]_{ij}$ [34]. Given $A = [A_{ij}] \in \mathbb{R}^{mp \times mp}$ with $A_{ij} \in \mathbb{R}^{m \times m}$, $\operatorname{vec}(A) \in \mathbb{R}^{m^2p^2}$ denotes the vectorization of A which stacks the columns of the matrix A, and

$$\operatorname{bvec}(\boldsymbol{A}) = [(\operatorname{vec}(\boldsymbol{A}_{11}))^{\top} (\operatorname{vec}(\boldsymbol{A}_{21}))^{\top} \cdots (\operatorname{vec}(\boldsymbol{A}_{p1}))^{\top} \\ (\operatorname{vec}(\boldsymbol{A}_{12}))^{\top} \cdots (\operatorname{vec}(\boldsymbol{A}_{p2}))^{\top} \cdots (\operatorname{vec}(\boldsymbol{A}_{pp}))^{\top}]^{\top}.$$

Let $S = \mathcal{E}_{\Delta} = \{\{k,\ell\} : \|\Delta^{(k\ell)}\|_F \neq 0\}$ where $\Delta = [\Delta^{(k\ell)}] \in \mathbb{R}^{mp \times mp}$ with $\Delta^{(k\ell)} \in \mathbb{R}^{m \times m}$ denoting the (k,l)th $m \times m$ submatrix of Δ . Then Δ_S denotes the submatrix of Δ with block rows and columns indexed by S, i.e., $\Delta_S = [\Delta^{(k\ell)}]_{(k,\ell) \in S}$. Suppose $\Gamma = A \boxtimes B$ given block partitioned matrices $A = [A_{ij}]$ and $B = [B_{k\ell}]$. For any two subsets T_1 and T_2 of $V \times V$, Γ_{T_1,T_2} denotes the submatrix of Γ with block rows and columns indexed by T_1 and T_2 , i.e., $\Gamma_{T_1,T_2} = [A_{j\ell} \otimes B_{kq}]_{(j,k) \in T_1,(\ell,q) \in T_2}$. Following [16], an operator $\mathcal{C}(\cdot)$ is used in Sec. IV. Consider $A \in \mathbb{R}^{mp \times mp}$ with (k,l)th $m \times m$ submatrix $A^{(k\ell)}$. Then $\mathcal{C}(\cdot)$ operates on A as

$$\begin{bmatrix} \boldsymbol{A}^{(11)} & \cdots & \boldsymbol{A}^{(1p)} \\ \vdots & \ddots & \vdots \\ \boldsymbol{A}^{(p1)} & \cdots & \boldsymbol{A}^{(pp)} \end{bmatrix} \xrightarrow{\boldsymbol{\mathcal{C}}(\cdot)} \begin{bmatrix} \|\boldsymbol{A}^{(11)}\|_F & \cdots & \|\boldsymbol{A}^{(1p)}\|_F \\ \vdots & \ddots & \vdots \\ \|\boldsymbol{A}^{(p1)}\|_F & \cdots & \|\boldsymbol{A}^{(pp)}\|_F \end{bmatrix}$$

with $\mathcal{C}(A^{(k\ell)}) = \|A^{(k\ell)}\|_F$ and $\mathcal{C}(A) \in \mathbb{R}^{p \times p}$. Now consider $A_1, A_2 \in \mathbb{R}^{mp \times mp}$ with (k,l)th $m \times m$ submatrices $A_1^{(k\ell)}$ and $A_2^{(k\ell)}$, respectively, and Tracy-Singh product $A_1 \boxtimes A_2 \in \mathbb{R}^{(mp)^2 \times (mp)^2}$. Then $\mathcal{C}(\cdot)$ operates on $A_1 \boxtimes A_2$ as $\mathcal{C}(A_1 \boxtimes A_2) \in \mathbb{R}^{p^2 \times p^2}$ with $\mathcal{C}(A_1^{(k_1\ell_1)} \otimes A_2^{(k_2\ell_2)}) = \|A_1^{(k_1\ell_1)} \otimes A_2^{(k_2\ell_2)}\|_F$ (= $\|A_1^{(k_1\ell_1)}\|_F \|A_2^{(k_2\ell_2)}\|_F$). That is, each $m^2 \times m^2$ submatrix $A_1^{(k_1\ell_1)} \otimes A_2^{(k_2\ell_2)}$ of $A_1 \boxtimes A_2$ is mapped into its Frobenius norm.

II. GROUP LASSO PENALIZED D-TRACE LOSS

Let $\boldsymbol{x} = [\boldsymbol{z}_{1x}^{\intercal} \ \boldsymbol{z}_{2x}^{\intercal} \ \cdots \ \boldsymbol{z}_{px}^{\intercal}]^{\intercal} \in \mathbb{R}^{mp}$ where $\boldsymbol{z}_{ix} \in \mathbb{R}^{m}$, $i \in [p]$, are zero-mean, jointly Gaussian. Similarly, let $\boldsymbol{y} = [\boldsymbol{z}_{1y}^{\intercal} \ \boldsymbol{z}_{2y}^{\intercal} \ \cdots \ \boldsymbol{z}_{py}^{\intercal}]^{\intercal} \in \mathbb{R}^{mp}$ where $\boldsymbol{z}_{iy} \in \mathbb{R}^{m}$, $i \in [p]$, are zero-mean, jointly Gaussian. Given i.i.d. samples $\boldsymbol{x}(t)$, $t = 1, 2, \cdots, n_{x}$, of \boldsymbol{x} , and similarly given i.i.d. samples $\boldsymbol{y}(t)$, $t = 1, 2, \cdots, n_{y}$, of $\boldsymbol{y} \in \mathbb{R}^{mp}$, form the sample covariance estimates

$$\hat{\boldsymbol{\Sigma}}_x = \frac{1}{n_x} \sum_{t=1}^{n_x} \boldsymbol{x}(t) \boldsymbol{x}^{\top}(t), \quad \hat{\boldsymbol{\Sigma}}_y = \frac{1}{n_y} \sum_{t=1}^{n_y} \boldsymbol{y}(t) \boldsymbol{y}^{\top}(t). \quad (2)$$

and denote their true values as $\Sigma_x^* = \Omega_x^{-*} (= (\Omega_x^*)^{-1})$ and $\Sigma_y^* = \Omega_y^{-*}$. Assume that $\{x(t)\}$ and $\{y(t)\}$ are mutually independent sequences. Assume Σ_x^* and Σ_y^* are positive definite. We wish to estimate $\Delta = \Omega_y^* - \Omega_x^*$ and graph $\mathcal{G}_{\Delta} = (V, \mathcal{E}_{\Delta})$, based on $\hat{\Sigma}_x$ and $\hat{\Sigma}_y$. Following the SA formulation of [8] (see also [28, Sec. 2.1]), we will use a convex D-trace loss function given by

$$L(\boldsymbol{\Delta}, \hat{\boldsymbol{\Sigma}}_x, \hat{\boldsymbol{\Sigma}}_y) = \frac{1}{2} \operatorname{tr}(\hat{\boldsymbol{\Sigma}}_x \boldsymbol{\Delta} \hat{\boldsymbol{\Sigma}}_y \boldsymbol{\Delta}^\top) - \operatorname{tr}(\boldsymbol{\Delta}(\hat{\boldsymbol{\Sigma}}_x - \hat{\boldsymbol{\Sigma}}_y)) \quad (3)$$

where D-trace refers to difference-in-trace loss function, a term coined in [29] in the context of graphical model estimation. The function $L(\Delta, \Sigma_x^*, \Sigma_y^*)$ is strictly convex in Δ (its Hessian w.r.t. $\operatorname{vec}(\Delta)$ is $\Sigma_y^* \otimes \Sigma_x^*$), and has a unique minimum at $\Delta^* = \Omega_y^* - \Omega_x^*$ [8], [28]. When we use sample covariances, we propose to estimate Δ by minimizing the group-lasso penalized loss function

$$L_{\lambda}(\boldsymbol{\Delta}, \hat{\boldsymbol{\Sigma}}_{x}, \hat{\boldsymbol{\Sigma}}_{y}) = L(\boldsymbol{\Delta}, \hat{\boldsymbol{\Sigma}}_{x}, \hat{\boldsymbol{\Sigma}}_{y}) + \lambda \sum_{k,\ell=1}^{p} \|\boldsymbol{\Delta}^{(k\ell)}\|_{F} \quad (4)$$

where $\lambda > 0$ is a tuning parameter and $\|\mathbf{\Delta}^{(k\ell)}\|_F$ promotes blockwise sparsity in $\mathbf{\Delta}$ [35]–[37] where, if we partition $\mathbf{\Delta}$

into $m \times m$ submatrices, $\mathbf{\Delta}^{(k\ell)}$ denotes its (k,ℓ) th submatrix, associated with edge $\{k,\ell\}$ of the differential graph $\mathcal{G}_{\Delta} = (V,\mathcal{E}_{\Delta})$.

For SA models (m=1), [28] has used the lasso-penalized loss function $L_J(\Delta) = L(\Delta, \hat{\Sigma}_x, \hat{\Sigma}_y) + \lambda \sum_{k,\ell=1}^p |\Delta_{k\ell}|$. The cost $L_J(\Delta)$ is optimized in [28] using a two-block ADMM approach which is known to be convergent. The resulting estimator $\hat{\Delta}$ that minimizes the above cost is not necessarily symmetric. To obtain a symmetric estimator for SA models, [8] proposes the lasso-penalized loss function $L_Y(\Delta) = \frac{1}{4} \text{tr}(\hat{\Sigma}_x \Delta \hat{\Sigma}_y \Delta^\top + \hat{\Sigma}_y \Delta \hat{\Sigma}_x \Delta^\top) - \text{tr}(\Delta(\hat{\Sigma}_x - \hat{\Sigma}_y)) + \lambda \sum_{k,\ell=1}^p |\Delta_{k\ell}|$. In [8], cost $L_Y(\Delta)$ is optimized using a three-block ADMM method which is not necessarily convergent.

Suppose

$$\hat{\mathbf{\Delta}} = \arg\min_{\mathbf{\Lambda}} L_{\lambda}(\mathbf{\Delta}, \hat{\mathbf{\Sigma}}_{x}, \hat{\mathbf{\Sigma}}_{y}). \tag{5}$$

Even though Δ is symmetric, $\hat{\Delta}$ is not. We can symmetrize it by setting $\hat{\Delta}_{sym} = \frac{1}{2}(\hat{\Delta} + \hat{\Delta}^{\top})$, after obtaining $\hat{\Delta}$. Then the differential graph edges are estimated as

$$\hat{\mathcal{E}}_{\Delta} = \left\{ \{k, \ell\} : \|\hat{\Delta}_{sym}^{(k\ell)}\|_F > 0 \right\}. \tag{6}$$

III. OPTIMIZATION

The objective function $L_{\lambda}(\Delta, \hat{\Sigma}_x, \hat{\Sigma}_y)$, given by (4), is strictly convex. Several existing approaches such as an alternating direction method of multipliers (ADMM) [38] or proximal gradient descent (PGD) methods [39], can be followed to minimize (4). Note that [8], [28] use ADMM while [9] uses a proximal gradient method, all for SA graphs. It is stated in [8, Sec. 2.2] that in their simulation example, ADMM approach yielded a slightly smaller value of the objective function compared to the PGD approach. In [10], [11], similar to [9], a proximal gradient method is used for an objective function similar to our (4). In this paper, motivated by [8], we will develop an ADMM method. In a simulation example (Sec. V-A) we compare our ADMM approach with ADMM and PGD approaches of [28] and [9], [11], respectively.

A. ADMM Approach

Similar to [28] (also [8]), we use an ADMM approach [38] with variable splitting. Using variable splitting, consider

$$\min_{\boldsymbol{\Delta}, \boldsymbol{W}} \left\{ L(\boldsymbol{\Delta}, \hat{\boldsymbol{\Sigma}}_x, \hat{\boldsymbol{\Sigma}}_y) + \lambda \sum_{k,\ell=1}^p \|\boldsymbol{W}^{(k\ell)}\|_F \right\}$$
(7)

subject to $\Delta = W$.

The scaled augmented Lagrangian for this problem is [38]

$$L_{\rho} = L(\boldsymbol{\Delta}, \hat{\boldsymbol{\Sigma}}_{x}, \hat{\boldsymbol{\Sigma}}_{y}) + \lambda \sum_{k,\ell=1}^{p} \|\boldsymbol{W}^{(k\ell)}\|_{F}$$
$$+ \frac{\rho}{2} \|\boldsymbol{\Delta} - \boldsymbol{W} + \boldsymbol{U}\|_{F}^{2}$$
(8)

where \boldsymbol{U} is the dual variable, and $\rho>0$ is the penalty parameter. Given the results $\boldsymbol{\Delta}^{(i)}, \boldsymbol{W}^{(i)}, \boldsymbol{U}^{(i)}$ of the ith iteration, in the (i+1)st iteration, an ADMM algorithm executes the following three updates:

(a)
$$\Delta^{(i+1)} \leftarrow \underset{L(\boldsymbol{\Delta}, \hat{\boldsymbol{\Sigma}}_x, \hat{\boldsymbol{\Sigma}}_y) + \frac{\rho}{2} \|\boldsymbol{\Delta} - \boldsymbol{W}^{(i)} + \boldsymbol{U}^{(i)}\|_F^2}{\arg \min_{\boldsymbol{\Delta}} L_a(\boldsymbol{\Delta})} := L_a(\boldsymbol{\Delta})$$

(b) $W^{(i+1)} \leftarrow \arg\min_{\boldsymbol{W}} L_b(\boldsymbol{W}), \ L_b(\boldsymbol{W}) := \lambda \sum_{k,\ell=1}^p \|\boldsymbol{W}^{(k\ell)}\|_F + \frac{\rho}{2} \|\boldsymbol{\Delta}^{(i+1)} - \boldsymbol{W} + \boldsymbol{U}^{(i)}\|_F^2.$ (c) $U^{(i+1)} \leftarrow U^{(i)} + (\boldsymbol{\Delta}^{(i+1)} - \boldsymbol{W}^{(i+1)}).$

Update (a): Differentiate $L_a(\Delta)$ w.r.t. Δ to obtain

$$\mathbf{0} = \frac{\partial L_a(\mathbf{\Delta})}{\partial \mathbf{\Delta}} = \hat{\mathbf{\Sigma}}_x \mathbf{\Delta} \hat{\mathbf{\Sigma}}_y - (\hat{\mathbf{\Sigma}}_x - \hat{\mathbf{\Sigma}}_y) + \rho(\mathbf{\Delta} - \mathbf{W} + \mathbf{U})$$
(9)

$$\Rightarrow (\hat{\Sigma}_y \otimes \hat{\Sigma}_x + \rho I) \text{vec}(\Delta) = \text{vec}(\hat{\Sigma}_x - \hat{\Sigma}_y + \rho (W - U))$$
(10)

Direct matrix inversion solution of (10) requires inversion of a $(mp)^2 \times (mp)^2$ matrix. A computationally cheaper solution is given in [8], [28], as follows. Carry out eigendecomposition of $\hat{\Sigma}_x$ and $\hat{\Sigma}_y$ as $\hat{\Sigma}_x = Q_x D_x Q_x^ op, \; Q_x Q_x^ op = I$ and $\hat{m{\Sigma}}_y = m{Q}_y m{Q}_y^ op, \; m{Q}_y m{Q}_y^ op = m{I}, \; ext{where} \; m{D}_x \; ext{and} \; m{D}_y \; ext{are}$ diagonal matrices of the respective eigenvalues. Then $\hat{\Delta}$ that minimizes $L_a(\Delta)$ is given by

$$\hat{\boldsymbol{\Delta}} = \boldsymbol{Q}_x \left[\boldsymbol{B} \circ \left[\boldsymbol{Q}_x^\top \left(\hat{\boldsymbol{\Sigma}}_x - \hat{\boldsymbol{\Sigma}}_y + \rho (\boldsymbol{W} - \boldsymbol{U}) \right) \boldsymbol{Q}_y \right] \right] \boldsymbol{Q}_y^\top \quad (11)$$

where the symbol \circ denotes the Hadamard product and $B \in$ $\mathbb{R}^{mp \times mp}$ organizes the diagonal of $(D_y \otimes D_x + \rho I)^{-1}$ in a matrix with $B_{jk} = 1/([D_x]_{jj}[D_y]_{kk} + \rho)$. Note that the eigendecomposition of $\hat{\Sigma}_x$ and $\hat{\Sigma}_y$ has to be done only once.

$$\boldsymbol{\Delta}^{(i+1)} = \boldsymbol{Q}_{x} \left[\boldsymbol{B} \circ \left[\boldsymbol{Q}_{x}^{\top} \left(\hat{\boldsymbol{\Sigma}}_{x} - \hat{\boldsymbol{\Sigma}}_{y} + \rho (\boldsymbol{W}^{(i)} - \boldsymbol{U}^{(i)}) \right) \boldsymbol{Q}_{y} \right] \right] \boldsymbol{Q}_{y}^{\top}$$
(12)

Update (b): Here we have the group lasso solution [35]–[37]

$$= \left(1 - \frac{(\lambda/\rho)}{\|(\boldsymbol{\Delta}^{(i+1)} + \boldsymbol{U}^{(i)})^{(k\ell)}\|_F}\right)_+ (\boldsymbol{\Delta}^{(i+1)} + \boldsymbol{U}^{(i)})^{(k\ell)}$$
(13)

where $(a)_{+} = \max(0, a)$.

A pseudocode for the ADMM algorithm, MA-ADMM, used in this paper is given in Algorithm 1 where we use the stopping (convergence) criterion following [38, Sec. 3.3.1] and varying penalty parameter ρ following [38, Sec. 3.4.1]. The stopping criterion is based on primal and dual residuals being small where, in our case, at (i + 1)st iteration, the primal residual is given by $\Delta^{(i+1)} - W^{(i+1)}$ and the dual residual by $\rho^{(i)}(\boldsymbol{W}^{(i+1)}-\boldsymbol{W}^{(i)})$. Convergence criterion is met when the norms of these residuals are below primary and dual tolerances τ_{pri} and τ_{dual} , respectively; see line 10 of Algorithm 1. In turn, τ_{pri} and τ_{dual} are chosen using an absolute and relative criterion as in line 10 of Algorithm 1 where τ_{abs} and τ_{rel} are user chosen absolute and relative tolerances, respectively. Line 10 of Algorithm 1 follows typical choices given in [38, Sec. 3.4.1]. For all numerical results presented later, we used $\rho_0 = 2$, $\mu = 10$, and $\tau_{abs} = \tau_{rel} = 10^{-4}$.

We will compare our approach with three other approaches in Sec. V-A. One of them is the single attribute (SA) based ADMM approach (see [8], [28]. A pseudocode of our implementation of this approach, SA-ADMM, is in Algorithm 2 which differs from in Algorithm 1 only in line 8 where we replace group lasso with elementwise lasso.

Algorithm 1 ADMM Algorithm MA-ADMM

Input: Data $\{x(t)\}_{t=1}^{n_x}$, $x \in \mathbb{R}^{mp}$, and $\{y(t)\}_{t=1}^{n_y}$, $y \in \mathbb{R}^{mp}$, regularization and penalty parameters λ and ρ_0 , tolerances τ_{abs} and τ_{rel} , variable penalty factor μ , maximum number of iterations i_{max} .

Output: estimated $\hat{\Delta}_{sym}$ and $\hat{\mathcal{E}}_{\Delta}$.

- 1: Calculate sample covariances $\hat{\boldsymbol{\Sigma}}_x = \frac{1}{n_x} \sum_{t=1}^{n_x} \boldsymbol{x}(t) \boldsymbol{x}^{\top}(t)$ and $\hat{\boldsymbol{\Sigma}}_y = \frac{1}{n_y} \sum_{t=1}^{n_y} \boldsymbol{y}(t) \boldsymbol{y}^\top(t)$. 2: Initialize: $\boldsymbol{\Delta}^{(0)} = \boldsymbol{U}^{(0)} = \boldsymbol{W}^{(0)} = \boldsymbol{0}$, where $\boldsymbol{\Delta}, \boldsymbol{U}, \boldsymbol{W} \in$
- $\mathbb{R}^{(mp)\times(mp)}, \, \rho^{(0)} = \rho_0.$
- 3: Eigendecompose $\hat{\Sigma}_x$ and $\hat{\Sigma}_y$ as $\hat{\Sigma}_x = Q_x D_x Q_x^{ op}$ and $\mathbf{\Sigma}_{y} = \mathbf{Q}_{y} \mathbf{D}_{y} \mathbf{Q}_{y}^{\perp}$.
- 4: converged = FALSE, i = 0
- 5: while converged = FALSE AND $i \leq i_{max}$, do
- Construct B $\mathbb{R}^{mp imes mp}$ with $oldsymbol{B}_{ik}$ $1/([\boldsymbol{D}_x]_{jj}[\boldsymbol{D}_y]_{kk} + \rho^{(i)}.$
- Set $\Delta^{(i+1)} = Q_x \Big[B \circ [Q_x^\top (\hat{\Sigma}_x \hat{\Sigma}_y + \rho(W^{(i)} \hat{\Sigma}_y)] \Big] \Big]$ $(U^{(i)})Q_y]Q_y^{ op}$.
- With $(a)_+ := \max(0,a)$, $\boldsymbol{A} = (\boldsymbol{\Delta}^{(i+1)} + \boldsymbol{U}^{(i)})^{(k\ell)}$ and $k, \ell \in [p]$, update $m \times m$ subblocks of W as

$$(\boldsymbol{W}^{(i+1)})^{(k\ell)} = \left(1 - \frac{(\lambda/\rho)}{\|\boldsymbol{A}\|_F}\right)_+ \boldsymbol{A}^{(k\ell)}.$$

- Dual update $U^{(i+1)} = U^{(i)} + (\Delta^{(i+1)} W^{(i+1)}).$
- Check convergence. Set tolerances

$$\tau_{pri} = mp \, \tau_{abs} + \tau_{rel} \, \max(\|\boldsymbol{\Delta}^{(i+1)}\|_F, \|\boldsymbol{W}^{(i+1)}\|_F)$$
$$\tau_{dual} = mp \, \tau_{abs} + \tau_{rel} \, \|\boldsymbol{U}^{(i+1)}\|_F / \rho^{(i)}.$$

Define $e_p = \| \boldsymbol{\Delta}^{(i+1)} - \boldsymbol{W}^{(i+1)} \|_F$ and $e_d = \rho^{(i)} \| \boldsymbol{W}^{(i+1)} - \boldsymbol{W}^{(i)} \|_F$. If $(e_p \leq \tau_{pri})$ AND $(e_d \leq$ τ_{dual}), set converged = TRUE .

- Update penalty parameter ρ : If $e_p > \mu e_d$, set $\rho^{(i+1)} =$ $2\rho^{(i)}$, else if $e_d > \mu e_p$, set $\rho^{(i+1)} = \rho^{(i)}/2$, otherwise $\rho^{(i+1)} = \rho^{(i)}$. We also need to set $U^{(i+1)} = U^{(i+1)}/2$ for $e_p > \mu e_d$ and $U^{(i+1)} = 2U^{(i+1)}$ for $e_d > \mu e_p$.
- $i \leftarrow i + 1$
- 13: end while
- 14: Set $\hat{\Delta}_{sym} = \frac{1}{2}(\boldsymbol{W} + \boldsymbol{W}^{\top})$. If $\|\hat{\Delta}_{sym}^{(jk)}\|_F > 0$, assign edge $\{j,k\} \in \hat{\mathcal{E}}_{\Delta}, \text{ else } \{j,k\} \not\in \hat{\mathcal{E}}_{\Delta}.$

Algorithm 2 ADMM Algorithm SA-ADMM

Input: As in Algorithm 1.

Output: Estimated $\hat{\Delta}_{sym}$ and $\hat{\mathcal{E}}_{\Delta}$.

- 1: Follow lines 1-7 of Algorithm 1
- 2: With $(a)_+ := \max(0,a), \ {m A} = ({m \Delta}^{(i+1)} + {m U}^{(i)})^{(k\ell)}$ and $k, \ell \in [mp]$, update $[\boldsymbol{W}]_{k\ell} \in \mathbb{R}$ as

$$[\boldsymbol{W}^{(i+1)}]_{k\ell} = \left(1 - \frac{(\lambda/\rho)}{|[\boldsymbol{A}]_{k\ell}|}\right)_{+} [\boldsymbol{A}]_{k\ell}.$$

3: Follow lines 9-14 of Algorithm 1

B. Proximal Gradient Descent Approach

It is a first-order method that is based on objective function values and gradient evaluations. A pseudocode of the PGD method of [11], MA-proximal, is in Algorithm 3, and that of [9], SA-proximal, is in Algorithm 4. Algorithm 4 differs from in Algorithm 3 only in line 8 where we replace group lasso with element-wise lasso. For all numerical results presented later, we used $\epsilon = 10^{-3}$ in line 7.

Algorithm 3 PGD Algorithm MA-PGD

Input: Data $\{x(t)\}_{t=1}^{n_x}$, $x \in \mathbb{R}^{mp}$, and $\{y(t)\}_{t=1}^{n_y}$, $y \in \mathbb{R}^{mp}$, tolerance ϵ , maximum number of iterations i_{max}

Output: Estimated $\hat{\Delta}_{sym}$ and $\hat{\mathcal{E}}_{\Delta}$.

- 1: Calculate sample covariances $\hat{\boldsymbol{\Sigma}}_x = \frac{1}{n_x} \sum_{t=1}^{n_x} \boldsymbol{x}(t) \boldsymbol{x}^\top(t)$ and $\hat{oldsymbol{\Sigma}}_y = rac{1}{n_y} \sum_{t=1}^{n_y} oldsymbol{y}(t) oldsymbol{y}^{ op}(t).$
- 2: Set $\eta = 1/(\phi_{max}(\hat{\Sigma}_x) \phi_{max}(\hat{\Sigma}_x))$. Initialize: $\Delta^{(0)} = 0$.
- 3: converged = FALSE, i = 0
- 4: while converged = FALSE AND $i \leq i_{max}$, do
- Set $\mathbf{A} = \mathbf{\Delta}^{(i)} \eta \Big(\hat{\mathbf{\Sigma}}_x \mathbf{\Delta}^{(i)} \hat{\mathbf{\Sigma}}_y (\hat{\mathbf{\Sigma}}_x \hat{\mathbf{\Sigma}}_y) \Big)$. For $k, \ell \in [p]$, update $m \times m$ subblocks as

$$(\boldsymbol{\Delta}^{(i+1)})^{(k\ell)} = \left(1 - \frac{\lambda \eta}{\|\boldsymbol{A}\|_F}\right)_+ \boldsymbol{A}^{(k\ell)} \,.$$

- $\text{If } \frac{L_{\lambda}(\boldsymbol{\Delta}^{(i+1)}, \hat{\boldsymbol{\Sigma}}_{x}, \hat{\boldsymbol{\Sigma}}_{y}) L_{\lambda}(\boldsymbol{\Delta}^{(i)}, \hat{\boldsymbol{\Sigma}}_{x}, \hat{\boldsymbol{\Sigma}}_{y})}{L_{\lambda}(\boldsymbol{\Delta}^{(i)}, \hat{\boldsymbol{\Sigma}}_{x}, \hat{\boldsymbol{\Sigma}}_{y})} \leq \epsilon \text{, set converged}$ = TRUE .
- $i \leftarrow i + 1$
- 9: end while
- 10: Set $\hat{\Delta}_{sym} = \frac{1}{2}(\Delta + \Delta^{\top})$. If $\|\hat{\Delta}_{sym}^{(jk)}\|_F > 0$, assign edge $\{j,k\} \in \hat{\mathcal{E}}_{\Delta}$, else $\{j,k\} \not\in \hat{\mathcal{E}}_{\Delta}$.

Algorithm 4 PGD Algorithm SA-PGD

Input: As in Algorithm 3 **Output:** Estimated $\hat{\Delta}$ and $\hat{\mathcal{E}}_{\Delta}$

- 1: Follow lines 1-5 of Algorithm 3
- 2: For $k, \ell \in [mp]$, update

$$[\boldsymbol{\Delta}^{(i+1)}]_{k\ell} = \left(1 - \frac{\lambda \eta}{|[\boldsymbol{A}]_{k\ell}|}\right)_{+} [\boldsymbol{A}]_{k\ell}.$$

3: Follow lines 7-10 of Algorithm 3

C. Computational Complexity

The computational complexity of ADMM and PGD methods has been discussed in [9] for SA differential graphs, and it is of the same order for MA graphs, because the difference lies only in lasso versus group lasso, i.e., element-wise softthresholding versus group-wise soft-thresholding. Noting that we have $mp \times mp$ precision matrices, by [9], the computational complexity of the ADMM approaches (our proposed and that of [8], [28]) is $\mathcal{O}((mp)^3)$ while that of the PGD methods of [9], [11] is either $\mathcal{O}((mp)^3)$ when as implemented in Algorithms 3 and 4, or $\mathcal{O}((n_x+n_y)(mp)^2)$ when an alternative implementation of the cost gradient in line 5 of Algorithms 3 and 4 is used (see [9, Sec. 2.2]). For $n_x + n_y \ge mp$, there is no advantage to this alternative approach.

D. Convergence of ADMM

The objective function (4), is strictly convex. It is also closed, proper and lower semi-continuous. Hence, for any fixed $\rho > 0$, the (two-block) ADMM algorithm is guaranteed to converge [38, Sec. 3.2], in the sense that we have primal residual convergence to 0, dual residual convergence to 0, and objective function convergence to the optimal value. For varying ρ , the convergence of ADMM has not been proven, but if we additionally impose $\rho^{(i)} = \rho^{(i_0)} > 0$ for $i \geq i_0$ for some i_0 , we have convergence [38, Sec. 3.4.1].

E. Model Selection

Following the lasso penalty work of [8] (who invokes [12]), we use the following criterion for selection of group lasso

$$BIC(\lambda) = (n_x + n_y) \|\hat{\mathbf{\Sigma}}_x \hat{\mathbf{\Delta}} \hat{\mathbf{\Sigma}}_y - (\hat{\mathbf{\Sigma}}_x - \hat{\mathbf{\Sigma}}_y)\|_F + \ln(n_x + n_y) |\hat{\mathbf{\Delta}}|_0$$
(14)

where $|A|_0$ denotes number of nonzero elements in A and $\hat{\Delta}$ obeys (5). Choose λ to minimize $BIC(\lambda)$. Following [8] we term it BIC (Bayesian information criterion) even though the cost function used is not negative log-likelihood although $\ln(n_x + n_y) |\hat{\Delta}|_0$ penalizes over-parametrization as in BIC. It is based on the fact that true Δ^* satisfies $\Sigma_x^* \Delta^* \Sigma_y^* - (\Sigma_x^* \Sigma_u^*$) = 0. Since (14) is not scale invariant, we scale both $\hat{\Sigma}_x^y$ and $\hat{\Sigma}_y$ (and $\hat{\Delta}$ commensurately) by $\bar{\Sigma}^{-1}$ where $\bar{\Sigma}=$ $\operatorname{diag}\{\hat{\Sigma}_x\}$ is a diagonal matrix of diagonal elements of $\hat{\Sigma}_x$.

In our simulations we search over $\lambda \in [\lambda_{\ell}, \lambda_u]$, where λ_{ℓ} and λ_u are selected via a heuristic as in [17]. Find the smallest λ , labeled λ_{sm} for which we get a no-edge model; then we set $\lambda_u = \lambda_{sm}/2$ and $\lambda_\ell = \lambda_u/10$.

IV. THEORETICAL ANALYSIS

Here we analyze the properties of $\hat{\Delta}$. Theorem 1 follows the approach(es) of [8], [16], [28], [29], [31] while Theorem 2 follows the general framework of [30]. The general method of [31] used in [8], [16], [28], [29] requires an irrepresentability condition (see (20)) which is not needed by the method of [30]. On the other hand, our Theorem 2 does not have a result like Theorem 1(ii), the oracle property, or support recovery Theorem 1(iv).

First some notation. Define the true differential edgeset

$$S = \mathcal{E}_{\Delta^*} = \{ \{k, \ell\} : \| \mathbf{\Delta}^{*(k\ell)} \|_F \neq 0 \}, \quad s = |S|.$$
 (15)

Define

$$\Gamma^* = \Sigma_y^* \boxtimes \Sigma_x^*, \quad \hat{\Gamma} = \hat{\Sigma}_y \boxtimes \hat{\Sigma}_x.$$
(16)

Also, recall the operator $\mathcal{C}(\cdot)$ defined in Sec. I-C. In the rest of this section, we allow p, s and λ to be a functions of sample size n, denoted as p_n , s_n and λ_n , respectively. Define

$$M = \max\{\|\mathcal{C}(\Sigma_x^*)\|_{\infty}, \|\mathcal{C}(\Sigma_y^*)\|_{\infty}\}, \tag{17}$$

$$M_{\Sigma} = \max\{\|\mathcal{C}(\mathbf{\Sigma}_x^*)\|_{1,\infty}, \|\mathcal{C}(\mathbf{\Sigma}_y^*)\|_{1,\infty}\}, \qquad (18)$$

$$\kappa_{\Gamma} = \| \mathcal{C}(\Gamma_{S,S}^*)^{-1} \|_{1,\infty}, \tag{19}$$

$$\alpha = 1 - \max_{e \in S^c} \| \mathcal{C}(\Gamma_{e,S}^*(\Gamma_{S,S}^*)^{-1}) \|_1,$$
 (20)

$$\bar{\sigma}_{xy} = \max_{i} \{ \max_{i} [\boldsymbol{\Sigma}_{x}^{*}]_{ii}, \max_{i} [\boldsymbol{\Sigma}_{y}^{*}]_{ii} \}$$
 (21)

$$C_0 = 40 \, m \, \bar{\sigma}_{xy} \sqrt{2(\tau + \ln(4m^2)/\ln(p_n))}$$
 (22)

where S and Γ^* have been defined in (15) and (16). In (20), we require $0 < \alpha < 1$, and the expression

$$\max_{e \in S^c} \|\mathcal{C}(\Gamma_{e,S}^*(\Gamma_{S,S}^*)^{-1})\|_1 \le 1 - \alpha$$

for some $\alpha \in (0,1)$ is called the *irrepresentability condition*. Similar conditions are also used in [8], [16], [28], [29], [31]. Let $\hat{\Delta}$ be as in (5).

Theorem 1. For the system model of Sec. II, under (15) and the irrepresentability condition (20) for some $\alpha \in (0, 1)$, if

$$\lambda_{n} = \max\left\{\frac{8}{\alpha}, \frac{3}{\alpha \bar{C}_{\alpha}} s_{n} \kappa_{\Gamma} M C_{M\kappa}\right\} C_{0} \sqrt{\frac{\ln(p_{n})}{n}}$$

$$n = \min(n_{x}, n_{y}) > \max\left\{\frac{1}{\min\{M^{2}, 1\}}, 81 M^{2} s_{n}^{2} \kappa_{\Gamma}^{2}, \frac{9s_{n}^{2}}{(\alpha \bar{C}_{\alpha})^{2}} (\kappa_{\Gamma} M C_{M\kappa})^{2}\right\} C_{0}^{2} \ln(p_{n})$$
(24)

where $\bar{C}_{\alpha}=\frac{1-\alpha}{2(2M+1)-2\alpha M}$ and $C_{M\kappa}=1.5\big(1+\kappa_{\Gamma}\min\{s_nM^2,M_{\Sigma}^2\}\big)$, then with probability $>1-2/p_n^{\tau-2}$, for any $\tau>2$, we have

(i)
$$\|\mathcal{C}(\hat{\Delta} - \Delta^*)\|_{\infty} \leq (C_{b1} + C_{b2})C_0\sqrt{\frac{\ln(p_n)}{n}}$$

where $C_{b1} = 3\kappa_{\Gamma} \max\left\{\frac{8}{\alpha}, \frac{3}{\alpha \bar{C}_{\alpha}} s_n \kappa_{\Gamma} M C_{M\kappa}\right\}$
 $C_{b2} = 9s_n \kappa_{\Gamma}^2 M^2$.

- (ii) $\hat{\Delta}_{S^c} = \mathbf{0}$.
- (iii) $\|\mathcal{C}(\hat{\Delta} \Delta^*)\|_F \leq \sqrt{s_n} \|\mathcal{C}(\hat{\Delta} \Delta^*)\|_{\infty}$.

(iv) Additionally, if
$$\min_{(k,\ell)\in S} \|(\mathbf{\Delta}^*)^{(k\ell)}\|_F \ge 2(C_{b1} + C_{b2})C_0\sqrt{\frac{\ln(p_n)}{n}}$$
, then $P(\mathcal{G}_{\hat{\Delta}} = \mathcal{G}_{\Delta^*}) > 1 - 2/p_n^{\tau-2}$ (support recovery)

The proof of Theorem 1 is given in Appendix A.

Now we present Theorem 2 that follows the general framework of [30]. Let $\hat{\Delta}$ be as in (5).

Theorem 2. For the system model of Sec. II, under (15), if

$$\lambda_n \ge (4 + 6Ms_n C_1) C_0 \sqrt{\frac{\ln(p_n)}{n}}$$

$$n = \min(n_x, n_y) > \max\left\{\frac{1}{M^2}, \left(\frac{96Ms_n}{\phi_{min}^*}\right)^2\right\} C_0^2 \ln(p_n)$$
(26)

where $C_1 = \max_{\{k,\ell\} \in V \times V} \|(\boldsymbol{\Delta}^*)^{(k\ell)}\|_F$ and $\phi_{min}^* = \phi_{min}(\boldsymbol{\Sigma}_x^*)\phi_{min}(\boldsymbol{\Sigma}_y^*)$, then with probability $> 1 - 2/p_n^{\tau-2}$, for any $\tau > 2$, we have

$$\|\hat{\Delta} - \Delta^*\|_F \le \frac{12\sqrt{s_n}}{\phi_{min}^*} (4 + 6Ms_n C_1) C_0 \sqrt{\frac{\ln(p_n)}{n}} \bullet (27)$$

The proof of Theorem 2 is given in Appendix B. Note that $\|\mathcal{C}(\hat{\Delta} - \Delta^*)\|_F = \|\hat{\Delta} - \Delta^*\|_F$ when comparing Theorems 1 and 2.

Remark 1: Convergence Rate. If M, M_{Σ} and κ_{Γ} stay bounded with increasing sample size n, we have $\|\mathcal{C}(\hat{\Delta} - \Delta^*)\|_F = \mathcal{O}_P(s_n^{1.5}\sqrt{\ln(p_n)/n})$. Therefore, for $\|\mathcal{C}(\hat{\Delta} - \Delta^*)\|_F \to 0$ as $n \to \infty$, we must have $s_n^{1.5}\sqrt{\ln(p_n)/n} \to 0$. The SA results in [8] need $s_n^{3.5}\sqrt{\ln(p_n)/n} \to 0$ when we take into account the dependence of various constants on s_n in [8]. Notice that M_{Σ} constraints covariances Σ_x^* and Σ_y^* which

can be dense even if Ω_x^* and Ω_y^* are sparse (they need not be sparse for differential estimation), making them possibly unbounded with increasing sample size n. In this case we use $\min\{s_nM^2,M_\Sigma^2\}=s_nM^2$ in $C_{M\kappa}$ and C_{b1} , with M bounded, leading to $\|\mathcal{C}(\hat{\Delta}-\Delta^*)\|_F=\mathcal{O}_P(s_n^{2.5}\sqrt{\ln(p_n)/n})$. On the other hand, in Theorem 2, we always have $\|\mathcal{C}(\hat{\Delta}-\Delta^*)\|_F=\mathcal{O}_P(s_n^{1.5}\sqrt{\ln(p_n)/n})$. \square

Remark 2: Our results assume Gaussian data. Theorems 1 and 2 continue to hold for more general sub-Gaussian distributions (Gaussian distribution is a sub-Gaussian distribution) except that the zeros in the precision matrix (see (1)), or in the difference of two precision matrices, no longer signify conditional independence (or change in conditional independence) of the random vectors associated with the respective nodes; they only imply zero partial correlation. We state Lemma 1 in Appendix A for sub-Gaussian distributions, following [31, Lemma 1]. Lemma 2 is then specialized to Gaussian distributions by setting the sub-Gaussian parameter $\sigma_{sg} = 1$. If $\sigma_{sg} \neq 1$, then the only changes required in Theorems 1 and 2 is a scaling of C_0 in (22) where one needs to replace the factor of 40 with $8(1 + 4\sigma_{sg}^2)$. \square

Remark 3: Theorems 1 and 2 assume a constant number of attributes m, with only p, s and λ allowed to be functions of sample size n, and our Remark 1 reflects this fact. In terms of m, the bounds in Theorem 1(i) and Theorem 2 are $\mathcal{O}(m^3)$, which follows from $M = \mathcal{O}(m)$ and $C_{M\kappa} = \mathcal{O}(m^2)$ in Theorem 1, and $M = \mathcal{O}(m)$, $C_1 = \mathcal{O}(m)$ and $C_0 = \mathcal{O}(m)$ in Theorem 2. Therefore, for large m, one would need much higher number of samples n, and it is not clear how to circumvent this bottleneck. A different class of models based on matrix-valued graphical modeling [40]-[42] is a potential solution. In a matrix graph set-up, one has matrix-valued observations Z, which in our context would require attributes (components of z_i) to be arranged along rows and nodes i along columns, with the covariance of vec(Z) having a Kronecker-product structure. This structure drastically reduces the number of unknowns from $\mathcal{O}((mp)^2)$ to $\mathcal{O}(m^2+p^2)$. Prior reported work ([40]-[42]) is on matrix graph estimation, with no reported work on differential matrix graphs.

V. NUMERICAL EXAMPLES

We now present numerical results for both synthetic and real data to illustrate the proposed approach. In synthetic data examples the ground truth is known and this allows for assessment of the efficacy of various approaches. In real data examples where the ground truth is unknown, our goal is visualization and exploration of the differential conditional dependency structures underlying the data.

A. Synthetic Data: Erdös-Rènyi and Barabási-Albert Graphs

We consider two types of graphs: Erdös-Rènyi (ER) graph and Barabási-Albert (BA) graph [43], [44]. The BA graphs are an example of scale-free graphs with power law degree distribution [43]. In the ER graph, p=100 nodes are connected to each other with probability $p_{er}=0.5$ and there are m=3 attributes per node whereas in the BA graph, we used p=100 and mean degree of 2 to generate a BA graph using the procedure given in [44] (MATLAB

function BAmodel.m from https://github.com/ShanLu1984/ Scale-Free-Network-Generation-and-Comparison). In the upper triangular Ω_x , we set $[\Omega_x^{(jk)}]_{st} = 0.5^{|s-t|}$ for j = k = $1, \dots, p, s, t = 1, \dots, m$. For $j \neq k$, if the two nodes are not connected in the graph (ER or BA), we have $\Omega^{(jk)} = 0$, and if nodes j and k are connected, then $[\Omega^{(jk)}]_{st}$ is uniformly distributed over $[-0.4, -0.1] \cup [0.1, 0.4]$. Then add lower triangular elements to make Ω_x a symmetric matrix. To generate Ω_{y} , we follow [8] and first generate a differential graph with $\Delta \in \mathbb{R}^{(mp) \times (mp)}$ as an ER graph (regardless of whether Ω_x is based on ER or BA model), with connection probability $p_{er} = 0.05$ (sparse): if nodes j and k are connected in the Ω_x model, then each of m^2 elements of $\Delta^{(jk)}$ is independently set to ± 0.9 with equal probabilities. Then $\Omega_y = \Omega_x + \Delta$. Finally add γI to Ω_y and to Ω_x and pick $\underline{\gamma}$ so that Ω_y and Ω_x are both positive definite. With $\Phi_x\Phi_x^{\top}=\Omega_x^{-1}$, we generate $x = \Phi w$ with $w \in \mathbb{R}^{mp}$ as zero-mean Gaussian, with identity covariance, and similarly for y. We generate $n = n_x = n_y$ i.i.d. observations for x and y, with m = 3, $p = 100, n \in \{100, 200, 300, 400, 800, 1200, 1600\}.$

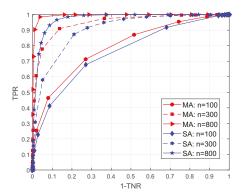


Fig. 1: ROC curves for ER graph based on ADMM approaches. TPR=true positive rate, TNR=true negative rate

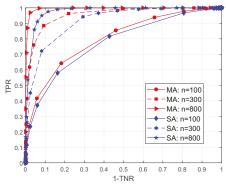


Fig. 2: ROC curves for BA graph based on ADMM approaches. TPR=true positive rate, TNR=true negative rate

Simulation results based on 100 runs are shown in Figs. 1-4. By changing the penalty parameter λ and determining the resulting edges, we calculated the true positive rate (TPR) and false positive rate 1-TNR (where TNR is the true negative rate) over 100 runs. The receiver operating characteristic (ROC) for ER graphs is shown in Fig. 1 for our MA-ADMM approach (labeled "MA") as well as for a SA-ADMM approach (labeled "SA"), based on [28], where we first estimate an

mp-node differential graph, and then use $\|\hat{\Delta}^{(k\ell)}\|_F \neq 0 \Leftrightarrow \{\{k,\ell\} \in \mathcal{E}_{\Delta}.$ It is seen from Fig. 1 that our MA-ADMM approach outperforms the SA-ADMM approach (that uses the same cost but element-wise lasso penalty instead of grouplasso penalty). Fig. 2 is the counterpart of Fig. 1 for BA graphs., and comments made regarding Fig. 1 apply here too.

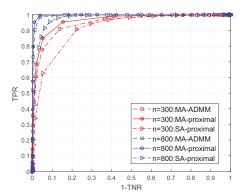


Fig. 3: ROC curves for ER graph based on ADMM as well as proximal approaches.

In Fig. 3 we compare ROC curves of our MA-ADMM approach with that for the MA-proximal and SA-proximal approaches of [11] and [9], respectively, for n = 300 and 800. It is seen that the MA-proximal approach outperforms our MA-ADMM approach, while both significantly outperform the SA-proximal approach (and the SA-ADMM approach whose ROC curves are in Fig. 1). In Table I, for the ER graph (n = 300, 800, p = 100, m = 3) we compare the four approaches (MA-ADMM, SA-ADMM, MA-proximal, SAproximal) in terms of the F_1 score, execution time (based on tic-toc functions in MATLAB), TPR and 1-TNR, for fixed penalty parameter λ selected from a grid of values (the same as for computing the ROC curves) to maximize the F_1 score averaged over 100 runs. All algorithms were run on a Window 10 Pro operating system with processor Intel(R) Core(TM) i7-10700 CPU @2.90 GHz with 32 GB RAM, using MATLAB R2023a. Notice that while the MA-proximal approach outperforms our MA-ADMM approach, it also takes more than twice the computation time for the MA-ADMM approach. Similarly, the SA-proximal approach takes more than twice the computation time for the SA-ADMM approach. The latter observation is consistent with the findings of [9].

In Table I, we also show results for n=3000 and 6000 for MA-ADMM and MA-proximal approaches in order to provide further empirical validation of the theoretical results stated in Theorem 1. Theorem 1 states that for sufficiently large sample size n, one can recover the differential graph structure exactly w.h.p. It is seen that the F_1 score approaches one with increasing n, implying graph support recovery, as claimed in Theorem 1(iv). Note that Theorem 1(iv) holds w.h.p., not with probability one, implying a nonzero probability of possibly inexact graph recovery, yielding an F1-score less than 1 for n=3000,6000. The sample sizes of n=300,800 are not large enough to yield an F_1 -score close to 1. There is a lower bound (24) on n for Theorem 1(iv) to hold w.h.p. The sample sizes of n=300,800 are apparently less than bound (which is not easily computable since it needs α and κ_{Γ}).

Approach	F_1 score $(\pm \sigma)$	timing (s) $(\pm \sigma)$	TPR $(\pm \sigma)$	1-TNR $(\pm \sigma)$
		n = 300		
MA-ADMM	0.6152 ± 0.0705	2.5044 ± 0.2939	0.6067 ± 0.1230	0.0184 ± 0.0080
MA-proximal	0.6686 ± 0.0639	6.6931 ± 0.3743	0.6845 ± 0.1253	0.0184 ± 0.0085
SA-ADMM	0.4549 ± 0.0332	0.1739 ± 0.0173	0.5795 ± 0.1132	0.0506 ± 0.0186
SA-proximal	0.4772 ± 0.0328	0.3517 ± 0.0195	0.6263 ± 0.1157	0.0524 ± 0.0193
n = 800				
MA-ADMM	0.8537 ± 0.0491	2.1911 ± 0.0639	0.9037 ± 0.0703	0.0111 ± 0.0041
MA-proximal	0.8898 ± 0.0408	4.8277 ± 0.2647	0.9526 ± 0.0537	0.0009 ± 0.0041
SA-ADMM	0.6336 ± 0.0055	0.1510 ± 0.0081	0.7468 ± 0.1017	0.0316 ± 0.0090
SA-proximal	0.6612 ± 0.0401	0.2533 ± 0.0143	0.7917 ± 0.0931	0.0314 ± 0.0090
n = 3000				
MA-ADMM	0.9795 ± 0.0152	2.1436 ± 0.0624	0.9928 ± 0.0018	0.0012 ± 0.0041
MA-proximal	0.9914 ± 0.0106	3.8312 ± 0.3948	0.9964 ± 0.0121	0.0007 ± 0.0007
n = 6000				
MA-ADMM	0.9914 ± 0.0087	1.9459 ± 0.0544	0.9997 ± 0.0013	0.0008 ± 0.0009
MA-proximal	0.9979 ± 0.0037	3.517 ± 0.1829	0.9998 ± 0.0011	0.0002 ± 0.0004

TABLE I: Comparisons among various approaches: Erdös-Rènyi graph, n = 300, 800, 3000, 6000, p = 100, m = 3. Tuning parameter λ picked to yield the highest F_1 score. Results based on 100 runs.

In Fig. 4 we show the results based on 100 runs for our approach when BIC parameter selection method (Sec. III-E) is applied in conjunction with the MA-ADMM approach. Here we show the TPR, 1-TNR and F_1 score values along with the $\pm \sigma$ error bars. The proposed approach works well both in terms of F_1 score and TPR vs 1-TNR.

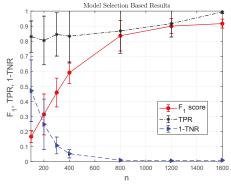


Fig. 4: BIC based results for ER graph: F_1 -scores, TPR and 1-TNR

B. Real Data: Beijing air-quality dataset [45]

Here we consider Beijing air-quality dataset [45], [46], downloaded from https://archive.ics.uci.edu/ml/datasets/ Beijing+Multi-Site+Air-Quality+Data. This data set includes hourly air pollutants data from 12 nationally-controlled airquality monitoring sites in the Beijing area from the Beijing Municipal Environmental Monitoring Center, and meteorological data in each air-quality site are matched with the nearest weather station from the China Meteorological Administration. The time period is from March 1st, 2013 to February 28th, 2017. The six air pollutants are PM_{2.5}, PM₁₀, SO₂, NO₂, CO, and O₃, and the meteorological data is comprised of five features: temperature, atmospheric pressure, dew point, wind speed, and rain; we did not use wind direction. Thus we have eleven features. We used data from 8 sites: 4 suburban/rural sites - Changping, Huairou, Shunyi, Dingling, and 4 urban area stations - Aotizhongxin, Dongsi, Guanyuan, Gucheng [46, Fig. 1]. The data are averaged over 24 hour period to yield daily averages. We used one year of daily data resulting in $n_x = n_y = 365$ days. The stations are used as attributes, with m = 8 for comparison between years 2013-14 and 2014-15, and m = 4 for comparison between suburban/rural sites and urban sites using 2013-14 year data.

We pre-process the data as follows. Given ith feature data $z_i(t) \in \mathbb{R}^m$, we transform it to $\bar{z}_i(t) = \ln(z_i(t)/z_i(t-1))$ and then detrend it (i.e., remove the best straight-line fit using the MATLAB function detrend). Finally, we scale the detrended scalar sequence to have a mean-square value of one over n_x or n_y samples. The logarithmic transformation and detrending of each feature sequence makes the sequence closer to (univariate) stationary and Gaussian, while scaling "balances" the possible wide variations in the scale of various feature measurements. All temperatures were converted from Celsius to Kelvin to avoid negative numbers, and if a value of a feature is zero (e.g., wind speed), we added a small positive number to it, so that the logarithmic transformation is well-defined.

Fig. 5 shows the estimated differential graphs when comparing daily-averaged data from 2013-14 (x-data) to that from 2014-15 (y-data), with air-quality and meteorological variables as p = 11 features measured at 8 monitoring sites (m=8). The objective is to visualize and explore differential conditional dependency relationships among the 11 variables, comparing one year to another, to investigate if pollution reduction measures have had any impact. Our intuition is that one does not expect such rapid changes within a short period of one year (see also [45]), therefore, our method should confirm our intuition. Figs. 5(a)-(b) show estimated $\|\hat{\Delta}^{(k\ell)}\|_F$ for various edges $\{k, \ell\}$, where it is unscaled in Fig. 5(a) but scaled in Fig. 5(b) so that the largest $\|\hat{\Delta}^{(k\ell)}\|_F$ (including $k=\ell$) is normalized to one. It is seen that differential graph weights are essentially zero (very sparse), implying that there are no yearto-year changes in the conditional dependency relationships among the 11 variables. This observation conforms to the findings of [45], [46]: no significant year-to-year changes. We also estimated the MA graphs for each year separately as $\|\hat{\Omega}_x^{(k\ell)}\|_F$ and $\|\hat{\Omega}_y^{(k\ell)}\|_F$ (shown in Figs. 5(c)-(d)), using the approach of [17], and based on the individual estimates, we computed the differential graph $\|\hat{\Omega}_y^{(k\ell)} - \hat{\Omega}_x^{(k\ell)}\|_F$ (shown in Figs. 5(e)-(f)). It is seen the separate-estimation approach

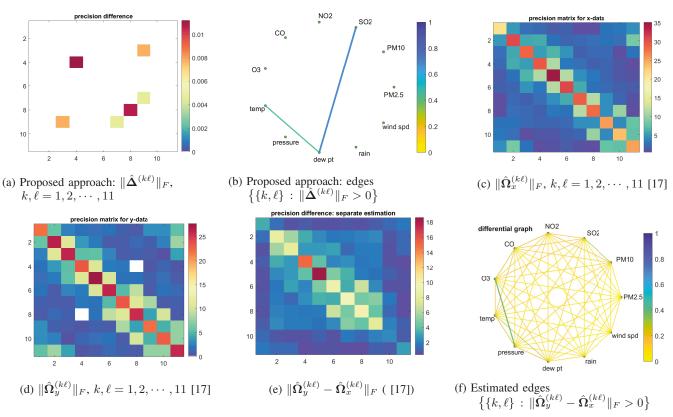


Fig. 5: Differential graphs comparing Beijing air-quality datasets [45] for years 2013-14 and 2014-15: 8 monitoring stations and 11 features ($m=8,\ p=11,\ n_x=n_y=365$). The features are numbered 1-11 beginning PM2.5 (PM_{2.5}) and moving counter-clockwise in Fig. 5(b).

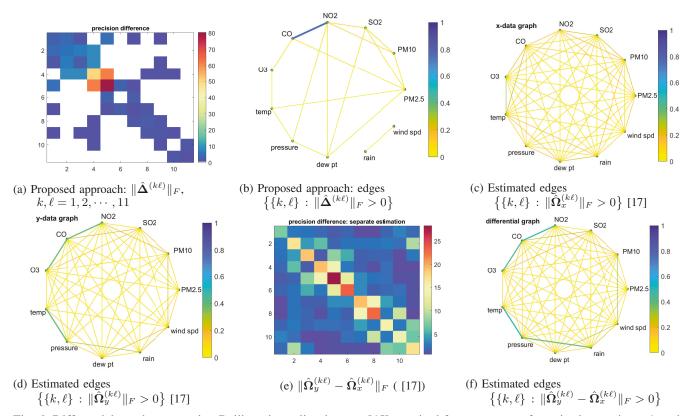


Fig. 6: Differential graphs comparing Beijing air-quality datasets [45] acquired from two sets of monitoring stations, 4 stations per set, year 2013-14: 4 monitoring stations and 11 features (m = 4, p = 11, $n_x = n_y = 365$).

does not yield a sparse differential graph, even though the two individual graphs in Figs. 5(c)-(d) are not all that different.

Fig. 6 shows the estimated differential graphs when comparing daily-averaged data over the period 2013-14, from four suburban/rural sites (x-data) to that from four urban sites (ydata), with air-quality and meteorological variables as p = 11features measured at two sets of 4 monitoring sites (m=4). The objective again is to visualize and explore differential conditional dependency relationships among the 11 variables, but in this case comparing one subregion to another. There are significant differences in meteorological conditions and pollutant sources, levels and mutual interactions, among suburban and urban areas [45], [46]. The suburban areas (located toward north) are less polluted than the urban areas (located toward south) [45], [46]. Automobile exhaust is the main cause of NO2 which is likely to undergo a chemical reaction with Ozone O₃, thereby, lowering its concentration [46]. Cold, dry air from the north reduces both dew point and PM_{2.5} particle concentration in suburban areas while southerly wind brings warmer and more humid air from the more polluted south that elevates the PM_{2.5} concentration [45]. The urban stations neighbor the south of Beijing which is heavily installed with iron, steel and cement industries in Hebei province [45]. Figs. 6(a)-(b) show estimated $\|\hat{\Delta}^{(k\ell)}\|_F$ for various edges $\{k,\ell\}$, where it is unscaled in Fig. 6(a) but scaled in Fig. 6(b) so that the largest $\|\hat{\Delta}^{(k\ell)}\|_F$ (including $k = \ell$) is normalized to one. It is seen that quite a few of the differential graph weights are significantly non-zero in Fig. 6(a), unlike that in Fig. 5(a), implying significant differences in the conditional dependency relationships among the 11 variables for suburban and urban areas. This observation conforms to the findings of [45], [46]. The comments made regarding Figs. 5(c)-(f) apply as well to Figs. 6(c)-(f).

VI. CONCLUSIONS

A group lasso penalized D-trace loss function approach for differential graph learning from multi-attribute data was presented. An ADMM algorithm was presented to optimize the convex objective function. Theoretical analysis establishing consistency of the estimator in high-dimensional settings was performed. We tested the proposed approach on synthetic as well as real data. In the synthetic data example, the multi-attribute approach is shown to outperform a single-attribute approach in correctly detecting the differential graph edges with ROC as the performance metric.

APPENDIX A TECHNICAL LEMMAS AND PROOF OF THEOREM 1

In this Appendix, we provide a proof of Theorem 1. A necessary and sufficient condition for minimization of convex $L_{\lambda}(\boldsymbol{\Delta}, \hat{\boldsymbol{\Sigma}}_{x}, \hat{\boldsymbol{\Sigma}}_{y})$ given by (4) w.r.t. $\boldsymbol{\Delta} \in \mathbb{R}^{mp \times mp}$ is that $\hat{\boldsymbol{\Delta}}$ minimizes (4) iff the zero matrix belongs to the sub-differential of $L_{\lambda}(\boldsymbol{\Delta}, \hat{\boldsymbol{\Sigma}}_{x}, \hat{\boldsymbol{\Sigma}}_{y})$. That is,

$$0 = \frac{\partial L(\boldsymbol{\Delta}, \hat{\boldsymbol{\Sigma}}_{x}, \hat{\boldsymbol{\Sigma}}_{y})}{\partial \boldsymbol{\Delta}} + \lambda \boldsymbol{Z}(\boldsymbol{\Delta}) \Big|_{\boldsymbol{\Delta} = \hat{\boldsymbol{\Delta}}}$$
$$= \hat{\boldsymbol{\Sigma}}_{x} \hat{\boldsymbol{\Delta}} \hat{\boldsymbol{\Sigma}}_{y} - (\hat{\boldsymbol{\Sigma}}_{x} - \hat{\boldsymbol{\Sigma}}_{y}) + \lambda \boldsymbol{Z}(\hat{\boldsymbol{\Delta}})$$
(28)

where $Z(\Delta) \in \partial \sum_{k,\ell=1}^p \|\Delta^{(k\ell)}\|_F \in \mathbb{R}^{mp \times mp}$, the sub-differential of group lasso penalty term, is given by [36], [37]

$$(Z(\Delta))^{(k\ell)} = \begin{cases} \frac{\Delta^{(k\ell)}}{\|\Delta^{(k\ell)}\|_F} & \text{if } \|\Delta^{(k\ell)}\|_F \neq 0 \\ V \in \mathbb{R}^{m \times m}, \|V\|_F \leq 1, & \text{if } \|\Delta^{(k\ell)}\|_F = 0 \end{cases} . (29)$$

In terms of $m \times m$ submatrices of Δ , $\hat{\Sigma}_x$, $\hat{\Sigma}_y$ and $Z(\Delta)$ corresponding to various graph edges, using $\operatorname{bvec}(ADB) = (B^{\top} \boxtimes A)\operatorname{bvec}(D)$ [33, Lemma 1], we may rewrite (28) as

$$(\hat{\mathbf{\Sigma}}_y \boxtimes \hat{\mathbf{\Sigma}}_x) \operatorname{bvec}(\hat{\mathbf{\Delta}}) - \operatorname{bvec}(\hat{\mathbf{\Sigma}}_x - \hat{\mathbf{\Sigma}}_y) + \lambda \operatorname{bvec}(\mathbf{Z}(\hat{\mathbf{\Delta}})) = \mathbf{0}$$
(30)

Then (30) can be rewritten as

$$\begin{bmatrix} \hat{\Gamma}_{S,S} & \hat{\Gamma}_{S,S^c} \\ \hat{\Gamma}_{S^c,S} & \hat{\Gamma}_{S^c,S^c} \end{bmatrix} \begin{bmatrix} \operatorname{bvec}(\hat{\Delta}_S) \\ \operatorname{bvec}(\hat{\Delta}_{S^c}) \end{bmatrix} - \begin{bmatrix} \operatorname{bvec}((\hat{\Sigma}_x - \hat{\Sigma}_y)_S) \\ \operatorname{bvec}((\hat{\Sigma}_x - \hat{\Sigma}_y)_{S^c}) \end{bmatrix} + \lambda \begin{bmatrix} \operatorname{bvec}(Z(\hat{\Delta}_S)) \\ \operatorname{bvec}(Z(\hat{\Delta}_S)) \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}.$$
(31)

The general approach of [31] (followed in [8], [16], [28], [29]) is to first solve the hypothetical constrained optimization problem with known edgeset S

$$\tilde{\Delta} = \arg\min_{\Delta: \Delta_{S^c} = \mathbf{0}} L_{\lambda}(\Delta, \hat{\Sigma}_x, \hat{\Sigma}_y)$$
 (32)

where S^c is the complement of S. Since, by construction, $\tilde{\Delta}_{S^c} = \mathbf{0}$, in this case (31) reduces to

$$\hat{\Gamma}_{S,S} \mathrm{bvec}(\tilde{\Delta}_S) - \mathrm{bvec}((\hat{\Sigma}_x - \hat{\Sigma}_y)_S) + \lambda \, \mathrm{bvec}(Z(\tilde{\Delta}_S)) = 0.$$
(33)

In the approach of [31], one investigates conditions under which the solution $\hat{\Delta}$ to (4) is the same as the solution $\tilde{\Delta}$ to (32). This is done by showing that $\hat{\Delta}$ satisfies (31). The choice $\hat{\Delta} = \tilde{\Delta}$ implies that $\hat{\Delta}_{S^c} = 0$ and (33) is true with $\tilde{\Delta}$ replaced with $\hat{\Delta}$. In order to satisfy (31), it remains to show that for any edge $e \in S^c$, we have strict feasibility

$$\|\hat{\mathbf{\Gamma}}_{e,S} \operatorname{bvec}(\tilde{\boldsymbol{\Delta}}_S) - \operatorname{bvec}((\hat{\boldsymbol{\Sigma}}_x - \hat{\boldsymbol{\Sigma}}_y)_e)\|_2 < \lambda,$$
 (34)

where for $a \in \mathbb{R}^q$, $||a||_2 = \sqrt{a^{\top}a}$. This requires a set of sufficient conditions stated in Theorem 1.

Lemma 1 follows from [31, Lemma 1]. It is stated for more general sub-Gaussian distributions as in [31, Lemma 1], but will be used later for Gaussian distributions, a subset of sub-Gaussian distributions.

Lemma 1. Suppose $\hat{\Sigma} = (1/n) \sum_{t=1}^n x(t) x^\top(t)$, given n i.i.d. samples $\{x(t)\}_{t=1}^n$ of zero-mean sub-Gaussian $x \in \mathbb{R}^{mp}$ with covariance Σ^* such that each component $x_i/\sqrt{\Sigma_{ii}^*}$ is sub-Gaussian with parameter σ_{sg} . Define $\sigma_{max} = \max_{1 \leq i \leq mp_n} \Sigma_{ii}^*$ and

$$\tilde{C}_0 = 8(1 + 4\sigma_{sg}^2)m\sigma_{max}\sqrt{2(\tau + \ln(4m^2)/\ln(p_n))}.$$
 (35)

Then

$$P\left(\|\mathcal{C}(\hat{\Sigma} - \Sigma^*)\|_{\infty} > \tilde{C}_0 \sqrt{\ln(p_n)/n}\right) \le 1/p_n^{\tau - 2}$$
 (36)

for any $\tau>2$ and $n>2m^2(\ln(4m^2)+\tau\ln(p_n))$. • Proof. By [31, Lemma 1], with $b=8(1+4\sigma_{sg}^2)$, we have

$$P(|[\hat{\Sigma} - \Sigma^*]_{ij}| > \delta) \le 4 \exp(-c_* n\delta^2)$$
 (37)

for any $\delta \in (0, b \, \sigma_{max})$ where $c_*^{-1} = 2b^2 \sigma_{max}^2$. For any edge $\{k, \ell\}$ of the MA graph, with m^2 edges $\{i, j\}$ of the corresponding SA graph associated with $\{k, \ell\}$, using the union bound, we have

$$P\left(|[\mathcal{C}(\hat{\Sigma} - \Sigma^*)]_{kl}| > \delta\right)$$

$$\leq P\left(\max_{\{i,j\}\in\{k,\ell\}} ([\hat{\Sigma} - \Sigma^*]_{ij})^2 > \frac{\delta^2}{m^2}\right)$$

$$\leq m^2 P\left(|[\hat{\Sigma} - \Sigma^*]_{ij}| > \frac{\delta}{m}\right) = 4m^2 \exp\left(-c_* n \frac{\delta^2}{m^2}\right).$$
(38)

Applying the union bound once more over all p_n^2 entries

$$P\Big(\|\mathcal{C}(\hat{\Sigma} - \Sigma^*)\|_{\infty} > \delta\Big) \le 4(mp_n)^2 \exp\Big(-c_* n \frac{\delta^2}{m^2}\Big) =: P_{tb}.$$
(39)

Choose $\delta = \tilde{C}_0 \sqrt{\ln(p_n)/n} = bm\sigma_{max} \sqrt{2\ln(4p_n^\tau m^2)/n}$. Then we have

$$P_{tb} = 4(mp_n)^2 \exp\left(\ln(4p_n^{\tau}m^2)^{-1}\right) = 1/p_n^{\tau-2} \tag{40}$$

provided $\delta \in (0,b\sigma_{max})$. Therefore, we need to have $\tilde{C}_0\sqrt{\ln(p_n)/n} < b\sigma_{max}$ requiring $n > 2m^2(\ln(4m^2) + \tau \ln(p_n))$. This completes the proof. \square

Using the union bound, Lemma 1 and Gaussian assumption, we have Lemma 2.

Lemma 2. Let $\hat{\Sigma}_x$ and $\hat{\Sigma}_y$ be as in (2), $\bar{\sigma}_{xy}$ as in (21), C_0 as in (22) and assume data are Gaussian. Define $n = \min(n_x, n_y)$ and

$$\mathcal{A} = \max \left\{ \| \mathcal{C}(\hat{\Sigma}_x - \Sigma_x^*) \|_{\infty}, \| \mathcal{C}(\hat{\Sigma}_y - \Sigma_y^*) \|_{\infty} \right\}. \tag{41}$$

Then for any $\tau > 2$ and $n > 2m^2 \ln(4m^2p_n^{\tau})$,

$$P(A > C_0 \sqrt{\ln(p_n)/n}) \le 2/p_n^{\tau-2} \quad \bullet$$
 (42)

Proof. For Gaussian distribution, the sub-Gaussian parameter σ_{sg} of Lemma 1 equals 1. Then $8(1+4\sigma_{sg}^2)=40$. Let $C_{0x}=40m(\max_i \Sigma_{x,ii}^*)\sqrt{2\big(\tau+\ln(4m^2)/\ln(p_n)\big)}$ and $C_{0y}=40m(\max_i \Sigma_{y,ii}^*)\sqrt{2\big(\tau+\ln(4m^2)/\ln(p_n)\big)}$ (where $\Sigma_{x,ii}^*=[\Sigma_x^*]_{ii}$, etc.). Then using Lemma 1 and union bound,

$$P\left(A > C_0 \sqrt{\ln(p_n)/n}\right)$$

$$\leq P\left(\|\mathcal{C}(\hat{\Sigma}_x - \Sigma_x^*)\|_{\infty} > C_0 \sqrt{\ln(p_n)/n}\right)$$

$$+ P\left(\|\mathcal{C}(\hat{\Sigma}_y - \Sigma_y^*)\|_{\infty} > C_0 \sqrt{\ln(p_n)/n}\right)$$

$$\leq 2/p_{\pi}^{\tau-2} \tag{43}$$

since $C_0 \ge C_{0x}$ and $C_0 \ge C_{0y}$. Recall (16)-(20) and define

 $\Delta_x = \hat{\Sigma}_x - \Sigma_x^*, \ \Delta_y = \hat{\Sigma}_y - \Sigma_y^*, \ \Delta_{\Gamma} = \hat{\Gamma} - \Gamma^*, \quad (44)$

$$\Delta_{\Sigma} = \Delta_x - \Delta_y, \ \epsilon_x = \|\mathcal{C}(\Delta_x)\|_{\infty}, \tag{45}$$

$$\epsilon_y = \|\mathcal{C}(\Delta_y)\|_{\infty}, \ \epsilon > \max\{\epsilon_x, \epsilon_y\}.$$
 (46)

Lemma 3. Assume that

$$\kappa_{\Gamma} < \frac{1}{3s_n(\epsilon^2 + 2M\epsilon)} \,. \tag{47}$$

Let $(\Gamma_{S,S}^{-*}$ denotes $(\Gamma_{S,S}^*)^{-1}$)

$$R(\boldsymbol{\Delta}_{\Gamma}) = \hat{\boldsymbol{\Gamma}}_{S,S}^{-1} - \boldsymbol{\Gamma}_{S,S}^{-*} + \boldsymbol{\Gamma}_{S,S}^{-*}(\boldsymbol{\Delta}_{\Gamma})_{S,S} \boldsymbol{\Gamma}_{S,S}^{-*}. \tag{48}$$

Then we have

$$\|\mathcal{C}(R(\Delta_{\Gamma}))\|_{\infty} \le \frac{3}{2} \kappa_{\Gamma}^3 s_n (\epsilon^2 + 2M\epsilon)^2, \tag{49}$$

$$\|\mathcal{C}(R(\Delta_{\Gamma}))\|_{1,\infty} \le \frac{3}{2}\kappa_{\Gamma}^3 s_n^2 (\epsilon^2 + 2M\epsilon)^2, \tag{50}$$

$$\|\mathcal{C}(\hat{\Gamma}_{S,S}^{-1} - \Gamma_{S,S}^{-*})\|_{\infty}$$

$$\leq \kappa_{\Gamma}^{2}(\epsilon^{2} + 2M\epsilon) \left(1 + 1.5 s_{n} \kappa_{\Gamma}(\epsilon^{2} + 2M\epsilon) \right), \tag{51}$$

$$\|\mathcal{C}(\hat{\Gamma}_{S,S}^{-1} - \Gamma_{S,S}^{-*})\|_{1,\infty} \le s_n \|\mathcal{C}(\hat{\Gamma}_{S,S}^{-1} - \Gamma_{S,S}^{-*})\|_{\infty}$$
 (52)

Proof. We have

$$\Delta_{\Gamma} = \hat{\Sigma}_{y} \boxtimes \hat{\Sigma}_{x} - \Sigma_{y}^{*} \boxtimes \Sigma_{x}^{*}
= \Delta_{y} \boxtimes \Delta_{x} + \Sigma_{y}^{*} \boxtimes \Delta_{x} + \Delta_{y} \boxtimes \Sigma_{x}^{*}.$$
(53)

By [16, Lemma 14],

$$\|\mathcal{C}(AB)\|_{1,\infty} \le \|\mathcal{C}(A)\|_{1,\infty} \|\mathcal{C}(B)\|_{1,\infty} \tag{54}$$

and by [16, Lemma 15],

$$\|\mathcal{C}(AB)\|_{\infty} \le \|\mathcal{C}(A)\|_{\infty} \|\mathcal{C}(B^{\top})\|_{1,\infty}. \tag{55}$$

Since $\|A \otimes B\|_F = \|A\|_F \|B\|_F$ and $A \boxtimes B = [A_{ij} \boxtimes B]_{ij} = [[A_{ij} \otimes B_{k\ell}]_{i\ell}]_{ij}$, we have

$$\|\mathcal{C}(A \boxtimes B)\|_{\infty} \le \|\mathcal{C}(A)\|_{\infty} \|\mathcal{C}(B)\|_{\infty}.$$
 (56)

From (17), (46), (53) and (56),

$$\|\mathcal{C}(\Delta_{\Gamma})\|_{\infty} \le \epsilon_x \epsilon_y + M \epsilon_x + M \epsilon_y < \epsilon^2 + 2M \epsilon$$
 (57)

and since $|S| = s_n$,

$$\|\mathcal{C}((\Delta_{\Gamma})_{S,S})\|_{1,\infty} \le s_n \|\mathcal{C}((\Delta_{\Gamma})_{S,S})\|_{\infty} \le s_n \|\mathcal{C}(\Delta_{\Gamma})\|_{\infty} < s_n(\epsilon^2 + 2M\epsilon).$$
 (58)

By assumption (47),

$$\kappa_{\Gamma} \| \mathcal{C}((\boldsymbol{\Delta}_{\Gamma})_{S,S}) \|_{1,\infty} = \| \mathcal{C}((\Gamma_{S,S}^*)^{-1}) \|_{1,\infty} \| \mathcal{C}((\boldsymbol{\Delta}_{\Gamma})_{S,S}) \|_{1,\infty}$$

$$< \frac{1}{3}.$$
(59)

By (59) we can invoke [31, Lemma 5] to have

$$R(\boldsymbol{\Delta}_{\Gamma}) = \boldsymbol{\Gamma}_{S,S}^{-*}(\boldsymbol{\Delta}_{\Gamma})_{S,S} \boldsymbol{\Gamma}_{S,S}^{-*}(\boldsymbol{\Delta}_{\Gamma})_{S,S} \boldsymbol{J} \boldsymbol{\Gamma}_{S,S}^{-*}$$
(60)

where $J=\sum_{k=0}^{\infty}(-1)^k \left(\Gamma_{S,S}^{-*}(\Delta_\Gamma)_{S,S}\right)^k$. Using (54), (55) and (60), we have

$$\|\mathcal{C}(R(\boldsymbol{\Delta}_{\Gamma}))\|_{\infty} \leq \|\mathcal{C}(\boldsymbol{\Gamma}_{S,S}^{-*}(\boldsymbol{\Delta}_{\Gamma})_{S,S})\|_{\infty} \times \|\mathcal{C}(\boldsymbol{\Gamma}_{S,S}^{-*}(\boldsymbol{\Delta}_{\Gamma})_{S,S}J\boldsymbol{\Gamma}_{S,S}^{-*})^{\top}\|_{1,\infty} \\ \leq \|\mathcal{C}(\boldsymbol{\Gamma}_{S,S}^{-*})\|_{1,\infty}^{3} \|\mathcal{C}((\boldsymbol{\Delta}_{\Gamma})_{S,S})\|_{\infty} \|\mathcal{C}((\boldsymbol{\Delta}_{\Gamma})_{S,S})\|_{1,\infty} \\ \times \|\mathcal{C}(\boldsymbol{J}^{\top})\|_{1,\infty}.$$
(61)

Now using (59),

$$\|\mathcal{C}(\boldsymbol{J}^{\top})\|_{1,\infty} \leq \sum_{k=0}^{\infty} \|\mathcal{C}(\boldsymbol{\Gamma}_{S,S}^{-*})\|_{1,\infty}^{k} \|\mathcal{C}((\boldsymbol{\Delta}_{\Gamma})_{S,S})\|_{1,\infty}^{k}$$

$$= \frac{1}{1 - \|\mathcal{C}(\boldsymbol{\Gamma}_{S,S}^{-*})\|_{1,\infty} \|\mathcal{C}((\boldsymbol{\Delta}_{\Gamma})_{S,S})\|_{1,\infty}}$$

$$= \frac{1}{1 - \kappa_{\Gamma} \|\mathcal{C}((\boldsymbol{\Delta}_{\Gamma})_{S,S})\|_{1,\infty}} \leq \frac{1}{1 - (1/3)} = \frac{3}{2}. \quad (62)$$

Using (57), (59), (61) and (62), we have

$$\|\mathcal{C}(R(\boldsymbol{\Delta}_{\Gamma}))\|_{\infty} \leq \frac{3}{2}\kappa_{\Gamma}^{3}s_{n}\|\mathcal{C}((\boldsymbol{\Delta}_{\Gamma})_{S,S})\|_{\infty}^{2}$$
$$<\frac{3}{2}\kappa_{\Gamma}^{3}s_{n}(\epsilon^{2}+2M\epsilon)^{2}. \tag{63}$$

This proves (49), from which (50) immediately follows. Using (19), (48), (54), (55) and (57) we have

$$\|\mathcal{C}(\hat{\Gamma}_{S,S}^{-1} - \Gamma_{S,S}^{-*})\|_{\infty} \leq \|\mathcal{C}(R(\Delta_{\Gamma}))\|_{\infty}$$

$$+ \|\mathcal{C}(\Gamma_{S,S}^{-*}(\Delta_{\Gamma})_{S,S}\Gamma_{S,S}^{-*})\|_{\infty}$$

$$\leq \|\mathcal{C}(R(\Delta_{\Gamma}))\|_{\infty} + \|\mathcal{C}(\Gamma_{S,S}^{-*})\|_{1,\infty}^{2} \|\mathcal{C}((\Delta_{\Gamma})_{S,S})\|_{\infty}$$

$$< \kappa_{\Gamma}^{2}(\epsilon^{2} + 2M\epsilon)(1 + 1.5s_{n}\kappa_{\Gamma}(\epsilon^{2} + 2M\epsilon)).$$
(64)

This proves (51). The claim (52) follows by noting that $|S| = s_n$. This completes the proof. \square

Lemma 4. Assume (47) and the following conditions:

$$0 < \alpha < 1$$
 where α is as in (20), (65)

$$\epsilon < \min \left\{ M, \frac{\alpha \lambda_n}{2(2-\alpha)} \right\},$$
(66)

$$\alpha C_{\alpha} \min\{\lambda_n, 1\} \ge 3s_n \epsilon M \kappa_{\Gamma} B_s \tag{67}$$

where

$$C_{\alpha} = \frac{\alpha \lambda_{n} + 2\epsilon \alpha - 4\epsilon}{2M\alpha \lambda_{n} + \alpha \lambda_{n} + 2\epsilon \alpha},$$

$$B_{s} = \left[1 + \kappa_{\Gamma} \left(3s_{n}\epsilon M + \min\{s_{n}M^{2}, M_{\Sigma}^{2}\}\right) \times \left(4.5s_{n}\epsilon M\kappa_{\Gamma} + 1\right)\right].$$
(68)

Then we have

(i) bvec($\hat{\Delta}_{S^c}$) = 0.

(ii)
$$\|\mathcal{C}(\hat{\Delta} - \hat{\Delta}^*)\|_{\infty} \le 2\lambda_n \kappa_{\Gamma} + 3s_n \epsilon M \kappa_{\Gamma}^2 (4.5s_n \epsilon M \kappa_{\Gamma} + 1)(2M + 2\lambda_n) \bullet$$

Proof. To establish part (i), we need to show that (34) is true. Let d denote the left-side of (34). It follow from (33) that

$$\operatorname{bvec}(\tilde{\boldsymbol{\Delta}}_{S}) = \hat{\boldsymbol{\Gamma}}_{S,S}^{-1} \left(\operatorname{bvec}((\hat{\boldsymbol{\Sigma}}_{x} - \hat{\boldsymbol{\Sigma}}_{y})_{S}) - \lambda \operatorname{bvec}(\boldsymbol{Z}(\tilde{\boldsymbol{\Delta}}_{S})) \right).$$
(70)

Substitute (70) in the left-side of (34) to yield

$$d = \|\hat{\boldsymbol{\Gamma}}_{e,S} \left[\hat{\boldsymbol{\Gamma}}_{S,S}^{-1} \left(\text{bvec}((\hat{\boldsymbol{\Sigma}}_x - \hat{\boldsymbol{\Sigma}}_y)_S) - \lambda \, \text{bvec}(\boldsymbol{Z}(\tilde{\boldsymbol{\Delta}}_S)) \right) \right] - \text{bvec}((\hat{\boldsymbol{\Sigma}}_x - \hat{\boldsymbol{\Sigma}}_y)_e) \|_2.$$
(71)

At the true values we have

$$\mathbf{0} = \frac{\partial L_{\lambda}(\boldsymbol{\Delta}, \boldsymbol{\Sigma}_{x}^{*}, \boldsymbol{\Sigma}_{y}^{*})}{\partial \boldsymbol{\Delta}} \bigg|_{\boldsymbol{\Delta} = \boldsymbol{\Delta}^{*}} = \boldsymbol{\Sigma}_{x}^{*} \boldsymbol{\Delta}^{*} \boldsymbol{\Sigma}_{y}^{*} - (\boldsymbol{\Sigma}_{x}^{*} - \boldsymbol{\Sigma}_{y}^{*})$$

implying

$$\Gamma^* \text{bvec}(\Delta^*) - \text{bvec}(\Sigma_x^* - \Sigma_y^*) = 0, \qquad (72)$$

which, noting that $(\Delta^*)_{S^c} = 0$, can be rewritten as (cf. (31))

$$\Gamma_{S,S}^* \operatorname{bvec}(\Delta_S^*) = \operatorname{bvec}(\Sigma_T^*)_S - \operatorname{bvec}(\Sigma_U^*)_S,$$
 (73)

$$\Gamma_{e,S}^* \text{bvec}(\Delta_S^*) = \text{bvec}(\Sigma_T^*)_e - \text{bvec}(\Sigma_U^*)_e$$
. (74)

Therefore, $(A^{-*} = (A^*)^{-1})$,

$$\begin{split} &\Gamma_{e,S}^* \Gamma_{S,S}^{-*} \big(\text{bvec}(\boldsymbol{\Sigma}_x^*)_S - \text{bvec}(\boldsymbol{\Sigma}_y^*)_S \big) \\ &- \text{bvec}(\boldsymbol{\Sigma}_x^*)_e + \text{bvec}(\boldsymbol{\Sigma}_y^*)_e = \mathbf{0} \,. \end{split} \tag{75}$$

Recalling (44) and using (75) in (71),

$$d = \|\hat{\boldsymbol{\Gamma}}_{e,S}\hat{\boldsymbol{\Gamma}}_{S,S}^{-1}\operatorname{bvec}((\boldsymbol{\Delta}_{\Sigma})_{S}) + (\hat{\boldsymbol{\Gamma}}_{e,S}\hat{\boldsymbol{\Gamma}}_{S,S}^{-1} - \boldsymbol{\Gamma}_{e,S}^{*}\boldsymbol{\Gamma}_{S,S}^{-*}) \left(\operatorname{bvec}(\boldsymbol{\Sigma}_{x}^{*})_{S} - \operatorname{bvec}(\boldsymbol{\Sigma}_{y}^{*})_{S}\right) - \lambda \hat{\boldsymbol{\Gamma}}_{e,S}\hat{\boldsymbol{\Gamma}}_{S,S}^{-1}\operatorname{bvec}(\boldsymbol{Z}(\tilde{\boldsymbol{\Delta}}_{S})) - \operatorname{bvec}((\boldsymbol{\Delta}_{\Sigma})_{e})\|_{2}.$$
 (76)

We now bound various terms in (76). Note that $\hat{\Gamma}_{e,S} \in \mathbb{R}^{m^2 \times (m^2 s_n)}$, $\hat{\Gamma}_{S,S}^{-1} \in \mathbb{R}^{(m^2 s_n) \times (m^2 s_n)}$, and $\operatorname{bvec}((\Delta_{\Sigma})_S) \in \mathbb{R}^{m^2 s_n}$ where $\Delta_{\Sigma} \in \mathbb{R}^{(mp_n) \times (mp_n)}$. Consider $A_{e,S} \in \mathbb{R}^{m^2 \times (m^2 s_n)}$. Then

$$\|\boldsymbol{A}_{e,S} \operatorname{bvec}((\boldsymbol{\Delta}_{\Sigma})_S)\|_2 = \|\sum_{f \in S} \boldsymbol{A}_{e,f} \operatorname{vec}((\boldsymbol{\Delta}_{\Sigma})_f)\|_2$$
 (77)

where edge $f \in S$, $\mathbf{A}_{e,f} \in \mathbb{R}^{m^2 \times m^2}$ and $(\mathbf{\Delta}_{\Sigma})_f \in \mathbb{R}^{m \times m}$. By the triangle inequality

$$\|\boldsymbol{A}_{e,S} \operatorname{bvec}((\boldsymbol{\Delta}_{\Sigma})_S)\|_2 \le \sum_{f \in S} \|\boldsymbol{A}_{e,f} \operatorname{vec}((\boldsymbol{\Delta}_{\Sigma})_f)\|_2.$$
 (78)

With $B_{i.}$ denoting the *i*th row of matrix B and using Cauchy-Schwartz inequality, we have

$$\|\mathbf{A}_{e,f} \operatorname{vec}((\mathbf{\Delta}_{\Sigma})_{f})\|_{2} = \left(\sum_{i=1}^{m^{2}} \left([\mathbf{A}_{e,f}]_{i.} \operatorname{vec}((\mathbf{\Delta}_{\Sigma})_{f}) \right)^{2} \right)^{1/2}$$

$$\leq \left(\sum_{i=1}^{m^{2}} \|[\mathbf{A}_{e,f}]_{i.}\|_{2}^{2} \|\operatorname{vec}((\mathbf{\Delta}_{\Sigma})_{f})\|_{2}^{2} \right)^{1/2}$$

$$= \|\operatorname{vec}((\mathbf{\Delta}_{\Sigma})_{f})\|_{2} \left(\sum_{i=1}^{m^{2}} \|[\mathbf{A}_{e,f}]_{i.}\|_{2}^{2} \right)^{1/2}$$

$$= \|\operatorname{vec}((\mathbf{\Delta}_{\Sigma})_{f})\|_{2} \|\mathbf{A}_{e,f}\|_{F}. \tag{79}$$

Therefore, using (78) and (79),

$$\|\boldsymbol{A}_{e,S}\operatorname{bvec}((\boldsymbol{\Delta}_{\Sigma})_{S})\|_{2} \leq \left(\sum_{f \in S} \|\boldsymbol{A}_{e,f}\|_{F}\right) \max_{g \in S} \|\operatorname{vec}((\boldsymbol{\Delta}_{\Sigma})_{g})\|_{2}$$
$$= \|\boldsymbol{\mathcal{C}}(\boldsymbol{A}_{e,S})\|_{1} \|\boldsymbol{\mathcal{C}}(\boldsymbol{\Delta}_{\Sigma})\|_{\infty}. \tag{80}$$

Using (80), we have

$$\begin{split} &\|\hat{\boldsymbol{\Gamma}}_{e,S}\hat{\boldsymbol{\Gamma}}_{S,S}^{-1}\text{bvec}((\boldsymbol{\Delta}_{\Sigma})_{S})\|_{2} \\ &\leq \|\boldsymbol{\mathcal{C}}(\hat{\boldsymbol{\Gamma}}_{e,S}\hat{\boldsymbol{\Gamma}}_{S,S}^{-1})\|_{1} \|\boldsymbol{\mathcal{C}}(\boldsymbol{\Delta}_{\Sigma})\|_{\infty}, \\ &\|(\hat{\boldsymbol{\Gamma}}_{e,S}\hat{\boldsymbol{\Gamma}}_{S,S}^{-1} - \boldsymbol{\Gamma}_{e,S}^{*}\boldsymbol{\Gamma}_{S,S}^{-*}) \big(\text{bvec}(\boldsymbol{\Sigma}_{x}^{*})_{S} - \text{bvec}(\boldsymbol{\Sigma}_{y}^{*})_{S}\big)\|_{2} \\ &\leq \|\boldsymbol{\mathcal{C}}(\hat{\boldsymbol{\Gamma}}_{e,S}\hat{\boldsymbol{\Gamma}}_{S,S}^{-1} - \boldsymbol{\Gamma}_{e,S}^{*}\boldsymbol{\Gamma}_{S,S}^{-*})\|_{1} \|\boldsymbol{\mathcal{C}}(\boldsymbol{\Sigma}_{x}^{*} - \boldsymbol{\Sigma}_{y}^{*})\|_{\infty}, \\ &\|\hat{\boldsymbol{\Gamma}}_{e,S}\hat{\boldsymbol{\Gamma}}_{S,S}^{-1}\text{bvec}(\boldsymbol{Z}(\tilde{\boldsymbol{\Delta}}_{S}))\|_{2} \\ &\leq \|\boldsymbol{\mathcal{C}}(\hat{\boldsymbol{\Gamma}}_{e,S}\hat{\boldsymbol{\Gamma}}_{S,S}^{-1})\|_{1} \|\boldsymbol{\mathcal{C}}(\boldsymbol{Z}(\tilde{\boldsymbol{\Delta}}_{S}))\|_{\infty}, \\ &\|\text{bvec}((\boldsymbol{\Delta}_{\Sigma})_{e})\|_{2} \leq \|\boldsymbol{\mathcal{C}}(\boldsymbol{\Delta}_{\Sigma})\|_{\infty}. \end{split} \tag{83}$$

By (29), (44) and (46)

$$\|\mathcal{C}(\Sigma_x^* - \Sigma_y^*)\|_{\infty} \le \|\mathcal{C}(\Sigma_x^*)\|_{\infty} + \|\mathcal{C}(\Sigma_y^*)\|_{\infty} \le 2M, \quad (85)$$

$$\|\mathcal{C}(Z(\tilde{\Delta}_S))\|_{\infty} \le 1,$$
 (86)

$$\|\mathcal{C}(\Delta_{\Sigma})\|_{\infty} \le \|\mathcal{C}(\Delta_x)\|_{\infty} + \|\mathcal{C}(\Delta_y)\|_{\infty} < 2\epsilon.$$
 (87)

Using (76) and (81)-87),

$$d < 2\epsilon \|\mathcal{C}(\hat{\Gamma}_{e,S}\hat{\Gamma}_{S,S}^{-1})\|_{1} + 2M\|\mathcal{C}(\hat{\Gamma}_{e,S}\hat{\Gamma}_{S,S}^{-1} - \Gamma_{e,S}^{*}\Gamma_{S,S}^{-*})\|_{1}$$
$$+ \lambda \|\mathcal{C}(\hat{\Gamma}_{e,S}\hat{\Gamma}_{S,S}^{-1})\|_{1} + 2\epsilon.$$
(88)

Therefore, $d < \lambda$ for any edge $e \in S^c$ if

$$U_{b1} := \max_{e \in S^{c}} 2M \| \mathcal{C}(\hat{\Gamma}_{e,S} \hat{\Gamma}_{S,S}^{-1} - \Gamma_{e,S}^{*} \Gamma_{S,S}^{-*}) \|_{1}$$

$$+ 2\epsilon (1 + \| \mathcal{C}(\hat{\Gamma}_{e,S} \hat{\Gamma}_{S,S}^{-1}) \|_{1}) \le \alpha \lambda_{n} (1 - C_{\alpha}), \quad (89)$$

$$U_{b2} := \max_{\alpha, \beta} \| \mathcal{C}(\hat{\Gamma}_{e,S} \hat{\Gamma}_{S,S}^{-1}) \|_{1} \le 1 - (1 - C_{\alpha})\alpha. \quad (90)$$

It remains to show that (89) and (90) are true under the assumptions of Lemma 4. Since

$$\hat{\Gamma}_{e,S}\hat{\Gamma}_{S,S}^{-1} - \Gamma_{e,S}^* \Gamma_{S,S}^{-*} = (\hat{\Gamma}_{e,S} - \Gamma_{e,S}^*) \Gamma_{S,S}^{-*}
+ \Gamma_{e,S}^* (\hat{\Gamma}_{S,S}^{-1} - \Gamma_{S,S}^{-*}) + (\hat{\Gamma}_{e,S} - \Gamma_{e,S}^*) (\hat{\Gamma}_{S,S}^{-1} - \Gamma_{S,S}^{-*}), (91)$$

we have

$$\begin{split} & \| \mathcal{C}(\hat{\Gamma}_{e,S} \hat{\Gamma}_{S,S}^{-1} - \Gamma_{e,S}^* \Gamma_{S,S}^{-*}) \|_{\infty} \\ & \leq \| \mathcal{C}(\hat{\Gamma}_{e,S} - \Gamma_{e,S}^*) \|_{\infty} \| \mathcal{C}(\Gamma_{S,S}^{-*}) \|_{1,\infty} \\ & + \| \mathcal{C}(\Gamma_{e,S}^*) \|_{\infty} \| \mathcal{C}(\hat{\Gamma}_{S,S}^{-1} - \Gamma_{S,S}^{-*}) \|_{1,\infty} \\ & + \| \mathcal{C}(\hat{\Gamma}_{e,S} - \Gamma_{e,S}^*) \|_{\infty} \| \mathcal{C}(\hat{\Gamma}_{S,S}^{-1} - \Gamma_{S,S}^{-*}) \|_{1,\infty} \,. \end{split}$$
(92)

With edge $e = \{i, k\} \in S^c$ and edge $f = \{j, \ell\} \in S$, consider

$$\hat{\Gamma}_{e,f} - \Gamma_{e,f}^* = \hat{\Gamma}_{ik,j\ell} - \Gamma_{ik,j\ell}^*
= \hat{\Sigma}_y^{(ij)} \otimes \hat{\Sigma}_x^{(k\ell)} - \Sigma_y^{*(ij)} \otimes \Sigma_x^{*(k\ell)}
= \Delta_y^{(ij)} \otimes \Delta_x^{(k\ell)} + \Sigma_y^{*(ij)} \otimes \Delta_x^{(k\ell)} + \Delta_y^{(ij)} \otimes \Sigma_x^{*(k\ell)}.$$
(93)

It then follows that

$$|\mathcal{C}(\hat{\Gamma}_{e,f} - \Gamma_{e,f}^*)| \leq ||\Delta_y^{(ij)}||_F ||\Delta_x^{(k\ell)}||_F + ||\Sigma_y^{*(ij)}||_F ||\Delta_x^{(k\ell)}||_F + ||\Delta_y^{(ij)}||_F ||\Sigma_x^{*(k\ell)}||_F \leq \epsilon_u \epsilon_x + M \epsilon_x + M \epsilon_y < \epsilon^2 + 2M \epsilon.$$
(94)

Hence

$$\|\mathcal{C}(\hat{\Gamma}_{e,S} - \Gamma_{e,S}^*)\|_{\infty} < \epsilon^2 + 2M\epsilon,$$
 (95)

$$\|\mathcal{C}(\hat{\Gamma}_{e,S} - \Gamma_{e,S}^*)\|_1 < s_n(\epsilon^2 + 2M\epsilon). \tag{96}$$

Since $\Gamma_{e,f}^* = \Sigma_y^{*(ij)} \otimes \Sigma_x^{*(k\ell)}$, we have $|\Gamma_{e,f}^*| \leq M^2$ and

$$\|\mathcal{C}(\Gamma_{e,S}^*)\|_{\infty} \le M^2, \quad \|\mathcal{C}(\Gamma_{e,S}^*)\|_1 \le s_n M^2.$$
 (97)

Alternatively, with $e = \{i, k\} \in S^c$ and $f = \{j, \ell\} \in S$,

$$\|\mathcal{C}(\Gamma_{e,S}^*)\|_1 = \sum_{f \in S} |\mathcal{C}(\Sigma_y^{*(ij)} \otimes \Sigma_x^{*(k\ell)})|$$

$$\leq \sum_{\{j,\ell\} \in S} \|\Sigma_y^{*(ij)}\|_F \|\Sigma_x^{*(k\ell)}\|_F$$

$$\leq (\sum_{j=1}^p \|\Sigma_y^{*(ij)}\|_F)(\sum_{\ell=1}^p \|\Sigma_x^{*(k\ell)}\|_F)$$

$$\leq \|\mathcal{C}(\Sigma_y^*)\|_{1,\infty} \|\mathcal{C}(\Sigma_x^*)\|_{1,\infty} \leq M_{\Sigma}^2.$$
 (98)

From (91)-(98) and Lemma 3, we have

$$\|\mathcal{C}(\hat{\Gamma}_{e,S}\hat{\Gamma}_{S,S}^{-1} - \Gamma_{e,S}^*\Gamma_{S,S}^{-*})\|_{1}$$

$$\leq \|\mathcal{C}(\hat{\Gamma}_{e,S} - \Gamma_{e,S}^*)\|_{1}\|\mathcal{C}(\Gamma_{S,S}^{-*})\|_{1,\infty}$$

$$+ \|\mathcal{C}(\Gamma_{e,S}^*)\|_{1}\|\mathcal{C}(\hat{\Gamma}_{S,S}^{-1} - \Gamma_{S,S}^{-*})\|_{1,\infty}$$

$$+ \|\mathcal{C}(\hat{\Gamma}_{e,S} - \Gamma_{e,S}^*)\|_{1}\|\mathcal{C}(\hat{\Gamma}_{S,S}^{-1} - \Gamma_{S,S}^{-*})\|_{1,\infty}$$

$$\leq s_{n}(\epsilon^{2} + 2M\epsilon)\kappa_{\Gamma} + \left[\min\{s_{n}M^{2}, M_{\Sigma}^{2}\} + s_{n}(\epsilon^{2} + 2M\epsilon)\right]$$

$$\times \left[s_{n}(\epsilon^{2} + 2M\epsilon)\kappa_{\Gamma}^{2}\right] (1 + 1.5s_{n}(\epsilon^{2} + 2M\epsilon)\kappa_{\Gamma})$$

$$\stackrel{\epsilon < M}{\leq} 3s_{n}\epsilon M\kappa_{\Gamma}B_{s} \leq \alpha C_{\alpha}\min\{\lambda_{n}, 1\}, \tag{99}$$

where B_s is as in (69) and we used $\epsilon < M$ to infer $\epsilon^2 + 2M\epsilon < 3M\epsilon$. Using the triangle inequality $|a| - |b| \le |a - b| \le |a| + |b|$, we have $\|\mathcal{C}(\hat{\Gamma}_{e,S}\hat{\Gamma}_{S,S}^{-1} - \Gamma_{e,S}^*\Gamma_{S,S}^{-*})\|_1 \ge \|\mathcal{C}(\hat{\Gamma}_{e,S}\hat{\Gamma}_{S,S}^{-1})\|_1 - \|\mathcal{C}(\Gamma_{e,S}^*\Gamma_{S,S}^{-*})\|_1$, which, using (65) and (99), leads to

$$\|\mathcal{C}(\hat{\Gamma}_{e,S}\hat{\Gamma}_{S,S}^{-1})\|_{1} \leq \|\mathcal{C}(\Gamma_{e,S}^{*}\Gamma_{S,S}^{-*})\|_{1} + \alpha C_{\alpha} \min\{\lambda_{n}, 1\}$$

$$\leq 1 - \alpha + \alpha C_{\alpha} \min\{\lambda_{n}, 1\} \leq 1 - (1 - C_{\alpha})\alpha. \quad (100)$$

This establishes (90). To show (89), using (99)-(100),

$$U_{b1} \leq 2M\alpha C_{\alpha} \min\{\lambda_{n}, 1\} + 2\epsilon (1 + 1 - (1 - C_{\alpha})\alpha)$$

$$\leq 2M\alpha C_{\alpha} \lambda_{n} + 2\epsilon (2 - (1 - C_{\alpha})\alpha) \stackrel{(68)}{=} \alpha \lambda_{n} (1 - C_{\alpha}).$$
(101)

This proves (89), and thus, part (i) of Lemma 4.

We now turn to the proof of Lemma 4(ii). Since $\hat{\Delta} = \tilde{\Delta}$, for any edge $\{k, \ell\} \in S$, we have

$$\|(\hat{\boldsymbol{\Delta}} - \boldsymbol{\Delta}^*)^{(k\ell)}\|_F = \|(\tilde{\boldsymbol{\Delta}} - \boldsymbol{\Delta}^*)^{(k\ell)}\|_F$$
$$= \|\operatorname{vec}(\tilde{\boldsymbol{\Delta}}^{(k\ell)}) - \operatorname{vec}((\boldsymbol{\Delta}^*)^{(k\ell)})\|_2. \tag{102}$$

Using (33) and (73)

$$\operatorname{bvec}((\tilde{\Delta} - \Delta^*)_S) = \hat{\Gamma}_{S,S}^{-1}\operatorname{bvec}((\Delta_{\Sigma})_S) + (\hat{\Gamma}_{S,S}^{-1} - \Gamma_{S,S}^{-*}) \times \operatorname{bvec}((\Sigma_x^* - \Sigma_y^*)_S) - \lambda_n \hat{\Gamma}_{S,S}^{-1}\operatorname{bvec}(Z(\tilde{\Delta}_S)).$$
(103)

Since
$$\hat{\Gamma}_{S,S}^{-1} = \hat{\Gamma}_{S,S}^{-1} - \Gamma_{S,S}^{-*} + \Gamma_{S,S}^{-*}$$
,

$$\|\mathcal{C}(\hat{\Gamma}_{S,S}^{-1})\|_{1,\infty} \le \|\mathcal{C}(\hat{\Gamma}_{S,S}^{-1} - \Gamma_{S,S}^{-*})\|_{1,\infty} + \|\mathcal{C}(\Gamma_{S,S}^{-*})\|_{1,\infty}.$$
(104)

By (103), for any edge $f = \{k, \ell\} \in S$, we have

$$\|\operatorname{vec}((\tilde{\Delta} - \Delta^*)^{(k\ell)})\|_{2} \leq \|(\hat{\Gamma}_{S,S}^{-1} - \Gamma_{S,S}^{-*})_{f,S} \times \operatorname{bvec}((\Delta_{\Sigma})_{S} + (\Sigma_{x}^{*} - \Sigma_{y}^{*})_{S} - \lambda_{n} Z(\tilde{\Delta}_{S}))\|_{2} + \|(\Gamma_{S,S}^{-*})_{f,S} \operatorname{bvec}((\Delta_{\Sigma})_{S} - \lambda_{n} Z(\tilde{\Delta}_{S}))\|_{2} \\ \leq \|\mathcal{C}(\hat{\Gamma}_{S,S}^{-1} - \Gamma_{S,S}^{-*})\|_{1,\infty} \left(\|\mathcal{C}(\Delta_{\Sigma})\|_{\infty} + \|\mathcal{C}(\Sigma_{x}^{*} - \Sigma_{y}^{*})\|_{\infty} + \lambda_{n}\right) + \|\mathcal{C}(\Gamma_{S,S}^{-*})\|_{1,\infty} \left(\|\mathcal{C}(\Delta_{\Sigma})\|_{\infty} + \lambda_{n}\right) \\ \leq s_{n} \kappa_{\Gamma}^{2}(\epsilon^{2} + 2M\epsilon)(1 + 1.5s_{n}(\epsilon^{2} + 2M\epsilon)\kappa_{\Gamma}) \\ \times (2\epsilon + 2M + \lambda_{n}) + \kappa_{\Gamma}(2\epsilon + \lambda_{n}) =: U_{b3}.$$
 (105)

By (66), for $0<\alpha<1$, we have $2\epsilon<\alpha\lambda_n/(2-\alpha)<\alpha\lambda_n<\lambda_n$. Therefore, $\kappa_\Gamma(2\epsilon+\lambda_n)<2\kappa_\Gamma\lambda_n$ and $2\epsilon+2M+\lambda_n<2M+2\lambda_n$. Since $\epsilon< M$ by (66), we also have $\epsilon^2+2M\epsilon<3M\epsilon$. Using these relations and (105), it follows that

$$U_{b3} \leq 3s_n \epsilon M \kappa_{\Gamma}^2 (1 + 4.5s_n \epsilon M \kappa_{\Gamma}) (2M + 2\lambda_n) + 2\lambda_n \kappa_{\Gamma}.$$

Finally,

$$\|\mathcal{C}(\hat{\boldsymbol{\Delta}} - \boldsymbol{\Delta}^*)\|_{\infty} = \max_{f = \{k,\ell\} \in S} \|\text{vec}((\tilde{\boldsymbol{\Delta}} - \boldsymbol{\Delta}^*)^{(k\ell)})\|_2,$$

proving the desired result. \Box

Proof of Theorem 1. Here we first show that under the sufficient conditions of Theorem 1, the assumptions of Lemmas 2-4 holds true. We pick $\epsilon = C_0 \sqrt{ln(p_n)/n}$, implying, by Lemma 2, that $\|\mathcal{C}(\hat{\Sigma}_x - \Sigma_x^*)\|_{\infty} \le \epsilon$ and $\|\mathcal{C}(\hat{\Sigma}_y - \Sigma_y^*)\|_{\infty} \le \epsilon$. with probability $\geq 1 - 2/p_n^{\tau-2}$, $\tau > 2$. We first show that with this choice of ϵ , condition (66) of Lemma 4 holds true. By the

choice of λ_n , we have $\lambda_n \geq 8\epsilon/\alpha$. Clearly, for $0 < \alpha < 1$, $(\alpha/8) < \alpha/(4(2-\alpha))$. Therefore

$$\epsilon \le \frac{\alpha \lambda_n}{8} < \frac{\alpha \lambda_n}{4(2-\alpha)} < \frac{\alpha \lambda_n}{2(2-\alpha)}$$
 (106)

By a choice of n in (24), we have $n > C_0^2 \ln(p_n)/M^2$. Hence, $\epsilon = C_0 \sqrt{\ln(p_n)/n} < M$. Thus, (66) of Lemma 4 holds true. Next we show that condition (47) of Lemma 3 holds. By a choice of n in (24), we have $n > 81C_0^2 \ln(p_n)M^2 s_n^2 \kappa_\Gamma^2$. Therefore,

$$\kappa_{\Gamma} < \frac{1}{C_0} \sqrt{\frac{n}{\ln(p_n)}} \times \frac{1}{9s_n M} = \frac{1}{9s_n M \epsilon} < \frac{1}{3s_n (\epsilon^2 + 2M\epsilon)}$$
(107)

since $\epsilon < M$. This proves (47) of Lemma 3 holds.

Now we show that (67) of Lemma 4 holds. Since $\epsilon < \alpha \lambda_n/(4(2-\alpha))$, by (68), we have

$$C_{\alpha} = \frac{\alpha \lambda_{n} + 2\epsilon \alpha - 4\epsilon}{2M\alpha \lambda_{n} + \alpha \lambda_{n} + 2\epsilon \alpha} > \frac{\alpha \lambda_{n} - 4\epsilon}{2M\alpha \lambda_{n} + \alpha \lambda_{n} + 2\epsilon \alpha}$$

$$> \frac{\alpha \lambda_{n} - \alpha \lambda_{n}/(2 - \alpha)}{\alpha \lambda_{n}(2M + 1) + 2\epsilon \alpha} = \frac{1 - \alpha}{(2 - \alpha)\left[2M + 1 + \frac{2\epsilon\alpha}{\lambda_{n}\alpha}\right]}$$

$$> \frac{1 - \alpha}{(2 - \alpha)(2M + 1) + \alpha} = \bar{C}_{\alpha}$$
(108)

where in the last inequality above, we used $\epsilon < \alpha \lambda_n/(2(2-\alpha))$ from (106). Consider the right-side $3s_n \epsilon M \kappa_{\Gamma} B_s$ of (67). From (107), $\kappa_{\Gamma} < 1/(9s_n M \epsilon)$. Therefore,

$$B_{s} < 1 + \kappa_{\Gamma} \left(\min\{s_{n}M^{2}, M_{\Sigma}^{2}\} + \frac{1}{3\kappa_{\Gamma}} \right) \left(\frac{1}{2} + 1 \right)$$

$$= 1.5 + 1.5\kappa_{\Gamma} \min\{s_{n}M^{2}, M_{\Sigma}^{2}\} = C_{M\kappa}. \tag{109}$$

By (23), $\lambda_n \geq 4.5\epsilon(\alpha \bar{C}_{\alpha})^{-1}s_nM\kappa_{\Gamma}(1+\kappa_{\Gamma}\min\{s_nM^2,M_{\Sigma}^2\})$. Hence, using (108), we have

$$3s_n \epsilon M \kappa_{\Gamma} B_s < \alpha \bar{C}_{\alpha} \lambda_n < \alpha C_{\alpha} \lambda_n$$
, (110)

proving part of (67) for our choice of λ_n . To show that we also have $3s_n \epsilon M \kappa_{\Gamma} B_s < \alpha \bar{C}_{\alpha}$, consider the choice of n in (24) given by

$$n > C_0^2 \ln(p_n) \frac{9s_n^2}{(\alpha \bar{C}_\alpha)^2} (\kappa_\Gamma M C_{M\kappa})^2.$$
 (111)

Then

$$\epsilon = C_0 \sqrt{\frac{\ln(p_n)}{n}} < \frac{\alpha \bar{C}_{\alpha}}{3s_n M \kappa_{\Gamma} C_{M\kappa}},$$
(112)

and from (110),

$$3s_n \epsilon M \kappa_{\Gamma} B_s < \alpha \bar{C}_{\alpha} < \alpha C_{\alpha} . \tag{113}$$

Thus, all assumptions of Lemma 4 hold true.

Therefore, Lemma 4(i) applies, proving Theorem 1(ii). By (107), $9s_n \epsilon M \kappa_{\Gamma} < 1$. Using this fact in Lemma 4(ii),

$$\|\mathcal{C}(\hat{\Delta} - \Delta^*)\|_{\infty} \le 2\lambda_n \kappa_{\Gamma} + 9s_n \epsilon M \kappa_{\Gamma}^2 (M + \lambda_n)$$

$$\le 3\lambda_n \kappa_{\Gamma} + 9s_n \epsilon M^2 \kappa_{\Gamma}^2. \tag{114}$$

Since, by (23), $3\lambda_n \kappa_\Gamma = C_0 \sqrt{\ln(p_n)/n} \, C_{b1}$ and we picked $\epsilon = C_0 \sqrt{\ln(p_n)/n}$, we have

$$\|\mathcal{C}(\hat{\Delta} - \Delta^*)\|_{\infty} \le (C_{b1} + C_{b2})C_0\sqrt{\ln(p_n)/n}$$

proving Theorem 1(i). To prove part (iii), since $\hat{\Delta}_{S^c} = \tilde{\Delta}_{S^c} = \Delta_{S^c}^* = 0$, we have

$$\|\mathcal{C}(\hat{\Delta} - \Delta^*)\|_F = \left(\sum_{\{k,\ell\} \in S} \|\hat{\Delta}^{(k\ell)} - (\Delta^*)^{(k\ell)}\|_F^2\right)^{1/2}$$

$$\leq \|\mathcal{C}(\hat{\Delta} - \Delta^*)\|_{\infty} \sqrt{s_n}. \tag{115}$$

Finally, to establish part (iv), note that parts (i)-(iii) hold with probability $> 1 - 2/p_n^{\tau-2}$ (with high probability (w.h.p.)). Recall that $\mathcal{G}_{\Delta} = (V, \mathcal{E}_{\Delta})$ denotes the MA differential graph with edgeset $\mathcal{E}_{\Delta} = \{\{k,\ell\} : \|\mathbf{\Delta}^{(k\ell)}\|_F > 0\}$. Let \mathcal{G}_{Δ^*} and $\hat{\mathcal{G}}_{\hat{\Delta}}$ denoted true and estimated graphs based on $\mathbf{\Delta}^*$ and $\hat{\mathbf{\Delta}}$, respectively. If $\min_{(k,\ell)\in S}\|(\mathbf{\Delta}^*)^{(k\ell)}\|_F \geq 2\|\mathcal{C}(\hat{\Delta} - \mathbf{\Delta}^*)\|_{\infty}$, then $\mathcal{C}(\hat{\Delta} - \mathbf{\Delta}^*)\|_{\infty} = \mathcal{C}((\hat{\Delta} - \mathbf{\Delta}^*)_S)\|_{\infty} \leq (1/2)\min_{(k,\ell)\in S}\|(\mathbf{\Delta}^*)^{(k\ell)}\|_F$, therefore, $\min_{(k,\ell)\in S}\|(\hat{\Delta}_S)^{(k\ell)}\|_F \geq (1/2)\min_{(k,\ell)\in S}\|(\mathbf{\Delta}^*)^{(k\ell)}\|_F > 0$, while $\hat{\mathbf{\Delta}}_{S^c} = \mathbf{0}$ w.h.p. \square

APPENDIX B

TECHNICAL LEMMAS AND PROOF OF THEOREM 2

In order to invoke [30], we first vectorize (3), using $\theta = \text{bvec}(\Delta) \in \mathbb{R}^{m^2p^2}$, as (cf. (30))

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \boldsymbol{\theta}^{\top} (\hat{\boldsymbol{\Sigma}}_{y} \boxtimes \hat{\boldsymbol{\Sigma}}_{x}) \boldsymbol{\theta} - \boldsymbol{\theta}^{\top} \text{bvec} (\hat{\boldsymbol{\Sigma}}_{x} - \hat{\boldsymbol{\Sigma}}_{y})$$
 (116)

where previous $L(\Delta, \hat{\Sigma}_x, \hat{\Sigma}_y)$ is now $\mathcal{L}(\theta)$. To include sparse-group penalty, recall that the submatrix $\Delta^{(k\ell)}$ of Δ corresponds to the edge $\{k,\ell\}$ of the MA graph. We denote its vectorized version as $\boldsymbol{\theta}_{Gt} \in \mathbb{R}^{m^2}$ (subscript G for grouped variables [30]) with index $t=1,2,\cdots,p^2$. Then $\boldsymbol{\theta}_{Gt}=\operatorname{vec}(\Delta^{(k\ell)})$ where $t=(k-1)p+\ell,\ \ell=t$ mod p, and $k=\lfloor t/p\rfloor+1$. Using this notation, the penalty $\lambda \sum_{k,\ell=1}^p \|\Delta^{(k\ell)}\|_F = \lambda \sum_{t=1}^{p^2} \|\boldsymbol{\theta}_{Gt}\|_2$. In the notation of [30], the regularization penalty

$$\mathcal{R}(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_{\bar{\mathcal{G}},2} := \sum_{t=1}^{p^2} \|\boldsymbol{\theta}_{Gt}\|_2$$
 (117)

where the index set $\{1, 2, \cdots, (mp)^2\}$ is partitioned into a set of $N_G = p^2$ disjoint groups $\bar{\mathcal{G}} = \{G_1, G_2, \cdots, G_{p^2}\}$. As shown in [30, Sec. 2.2], $\mathcal{R}(\theta)$ is a norm. The counterpart to penalized $L_{\lambda}(\boldsymbol{\Delta}, \hat{\boldsymbol{\Sigma}}_x, \hat{\boldsymbol{\Sigma}}_y)$ of (4) is (we denote λ by λ_n , as in Appendix A)

$$\mathcal{L}_{\lambda}(\boldsymbol{\theta}) = \mathcal{L}(\boldsymbol{\theta}) + \lambda_n \mathcal{R}(\boldsymbol{\theta}). \tag{118}$$

As discussed in [30, Sec. 2.2], w.r.t. the usual Euclidean inner product $\langle \boldsymbol{u}, \boldsymbol{v} \rangle = \boldsymbol{u}^{\top} \boldsymbol{v}$ for $\boldsymbol{u}, \boldsymbol{v} \in \mathbb{R}^{m^2 p^2}$ and given any subset $S_{\bar{\mathcal{G}}} \subseteq \{1, 2, \cdots, N_G\}$ of group indices, define the subspace

$$\mathcal{M} = \{ \boldsymbol{\theta} \in \mathbb{R}^{m^2 p^2} \, | \, \boldsymbol{\theta}_{Gt} = \mathbf{0} \text{ for all } t \notin S_{\bar{G}} \}$$
 (119)

and its orthogonal complement

$$\mathcal{M}^{\perp} = \{ \boldsymbol{\theta} \in \mathbb{R}^{m^2 p^2} \, | \, \boldsymbol{\theta}_{Gt} = \mathbf{0} \text{ for all } t \in S_{\bar{\mathcal{G}}} \}.$$
 (120)

The chosen $\mathcal{R}(\boldsymbol{\theta})$ is decomposable w.r.t. $(\mathcal{M}, \mathcal{M}^{\perp})$ since $\mathcal{R}(\boldsymbol{\theta}^{(1)} + \boldsymbol{\theta}^{(2)}) = \mathcal{R}(\boldsymbol{\theta}^{(1)}) + \mathcal{R}(\boldsymbol{\theta}^{(2)})$ for any $\boldsymbol{\theta}^{(1)} \in \mathcal{M}$ and $\boldsymbol{\theta}^{(2)} \in \mathcal{M}^{\perp}$ [30, Sec. 2.2, Example 2].

In order to invoke [30], we need the dual norm \mathcal{R}^* of regularizer \mathcal{R} w.r.t. the inner product $\langle u, v \rangle = u^\top v$. It is given by [30, Eqns. (14)-(15)]

$$\mathcal{R}^*(v) = \sup_{\mathcal{R}(u) < 1} \langle u, v \rangle = \max_{t=1, 2, \dots p^2} \|u_{Gt}\|_2.$$
 (121)

We also need the subspace compatibility index [30], defined as

$$\Psi(\mathcal{M}) = \sup_{\boldsymbol{u} \in \mathcal{M} \setminus \{0\}} \frac{\mathcal{R}(\boldsymbol{u})}{\|\boldsymbol{u}\|_2}.$$
 (122)

For group lasso penalty, $\Psi(\mathcal{M}) = \sqrt{s_n}$ [30, Sec. 9.2 (Supplementary)], where $s_n = |S_{\bar{\mathcal{G}}}| = \text{number of edges in the true MA differential graph. We need to establish a restricted strong convexity condition [30] on <math>\mathcal{L}(\theta)$. With $\theta^* = \text{bvec}(\Delta^*)$ denoting the true value, and $\theta = \theta^* + \tilde{\theta}$, define

$$\delta \mathcal{L}(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) = \mathcal{L}(\boldsymbol{\theta}^* + \tilde{\boldsymbol{\theta}}) - \mathcal{L}(\boldsymbol{\theta}^*) - \langle \nabla \mathcal{L}(\boldsymbol{\theta}^*), \tilde{\boldsymbol{\theta}} \rangle$$
 (123)

where the gradient $\nabla \mathcal{L}(\boldsymbol{\theta}^*)$ at $\boldsymbol{\theta} = \boldsymbol{\theta}^*$ is

$$\nabla \mathcal{L}(\boldsymbol{\theta}^*) = (\hat{\boldsymbol{\Sigma}}_y \boxtimes \hat{\boldsymbol{\Sigma}}_x) \boldsymbol{\theta}^* - \operatorname{bvec}(\hat{\boldsymbol{\Sigma}}_x - \hat{\boldsymbol{\Sigma}}_y). \tag{124}$$

Hence (123) simplifies to

$$\delta \mathcal{L}(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) = \tilde{\boldsymbol{\theta}}^\top (\hat{\boldsymbol{\Sigma}}_u \boxtimes \hat{\boldsymbol{\Sigma}}_x) \tilde{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}^\top \hat{\boldsymbol{\Gamma}} \tilde{\boldsymbol{\theta}}, \qquad (125)$$

which may be rewritten as

$$\delta \mathcal{L}(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) = \tilde{\boldsymbol{\theta}}^{\top} \boldsymbol{\Gamma}^* \tilde{\boldsymbol{\theta}} + \tilde{\boldsymbol{\theta}}^{\top} (\hat{\boldsymbol{\Gamma}} - \boldsymbol{\Gamma}^*) \tilde{\boldsymbol{\theta}}. \tag{126}$$

By the sparsity assumption, $\theta^* = \theta_{\mathcal{M}}^*$, hence, $\theta_{\mathcal{M}^{\perp}}^* = 0$, where $\theta_{\mathcal{M}}$ and $\theta_{\mathcal{M}^{\perp}}$ denote projection of θ on subspaces \mathcal{M} and \mathcal{M}^{\perp} , respectively.

Similar to (5), suppose

$$\hat{\boldsymbol{\theta}} = \arg\min_{\boldsymbol{\theta}} \left\{ \mathcal{L}(\boldsymbol{\theta}) + \lambda_n \mathcal{R}(\boldsymbol{\theta}) \right\}, \tag{127}$$

and we consider (123) and (125) with $\hat{\theta} = \theta^* + \tilde{\theta}$. Then

$$\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^* = \hat{\boldsymbol{\theta}}_{\mathcal{M}} - \boldsymbol{\theta}^* + \hat{\boldsymbol{\theta}}_{\mathcal{M}^{\perp}} = \tilde{\boldsymbol{\theta}}_{\mathcal{M}} + \tilde{\boldsymbol{\theta}}_{\mathcal{M}^{\perp}}. \tag{128}$$

By [30, Lemma 1],

$$\mathcal{R}(\tilde{\boldsymbol{\theta}}_{\mathcal{M}^{\perp}}) \leq 3\mathcal{R}(\tilde{\boldsymbol{\theta}}_{\mathcal{M}}) + 4\mathcal{R}(\boldsymbol{\theta}_{\mathcal{M}^{\perp}}^{*}), \qquad (129)$$

if

$$\lambda_n \ge 2\mathcal{R}^*(\nabla \mathcal{L}(\boldsymbol{\theta}^*)). \tag{130}$$

Since in our case $\boldsymbol{\theta}_{\mathcal{M}^{\perp}}^{*}=\mathbf{0}$, we have $\mathcal{R}(\boldsymbol{\theta}_{\mathcal{M}^{\perp}}^{*})=0$. **Lemma 5**. Under (15) and using the notation of Appendix A,

$$\mathcal{R}^*(\nabla \mathcal{L}(\boldsymbol{\theta}^*)) \leq (\epsilon^2 + 2M\epsilon) s_n \max_{t=1,\dots,n} \|\boldsymbol{\theta}_{Gt}^*\|_2 + 2\epsilon \bullet$$

Proof. Using (44), (72) and (124), we have

$$\nabla \mathcal{L}(\boldsymbol{\theta}^*) = \boldsymbol{\Delta}_{\Gamma} \boldsymbol{\theta}^* + \operatorname{bvec}(\boldsymbol{\Delta}_u) - \operatorname{bvec}(\boldsymbol{\Delta}_x). \tag{131}$$

Expressing it group-wise, with groups t and t_1 corresponding to edges $\{j, k\}$ and $\{\ell, q\}$, respectively,

$$(\nabla \mathcal{L}(\boldsymbol{\theta}^*))_{Gt_1} = \sum_{t=1}^{p^2} (\boldsymbol{\Delta}_{\Gamma})_{Gt_1,Gt} \boldsymbol{\theta}_{Gt}^* + \operatorname{bvec}(\boldsymbol{\Delta}_y)_{Gt_1} - \operatorname{bvec}(\boldsymbol{\Delta}_x)_{Gt_1}.$$
(132)

Therefore, by the Cauchy-Schwartz inequality, and using (45), (46) and (57), we have

$$\|(\nabla \mathcal{L}(\boldsymbol{\theta}^{*}))_{Gt_{1}}\|_{2} \leq \sum_{t=1}^{p^{2}} \|(\boldsymbol{\Delta}_{\Gamma})_{Gt_{1},Gt}\|_{F} \|\boldsymbol{\theta}_{Gt}^{*}\|_{2}$$

$$+ \|\operatorname{bvec}(\boldsymbol{\Delta}_{y})_{Gt_{1}}\|_{2} + \|\operatorname{bvec}(\boldsymbol{\Delta}_{x})_{Gt_{1}}\|_{2}$$

$$\leq \|\boldsymbol{\mathcal{C}}(\boldsymbol{\Delta}_{\Gamma})\|_{\infty} \sum_{t=1}^{p^{2}} \|\boldsymbol{\theta}_{Gt}^{*}\|_{2} + \|\boldsymbol{\Delta}_{y}^{(\ell q)}\|_{\infty} + \|\boldsymbol{\Delta}_{y}^{(\ell q)}\|_{\infty}$$

$$\leq (\epsilon^{2} + 2M\epsilon) s_{n} \max_{t=1,\dots,p^{2}} \|\boldsymbol{\theta}_{Gt}^{*}\|_{2} + \epsilon + \epsilon. \tag{133}$$

By (121) and (133) we have the desired result. □ **Lemma 6.** Under (15) and the notation of Appendix A,

$$\delta \mathcal{L}(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \ge \kappa_{\mathcal{L}} \|\tilde{\boldsymbol{\theta}}\|_2^2 \tag{134}$$

where $\kappa_{\mathcal{L}} = \frac{1}{2} \phi_{min}^* - 8s_n(\epsilon^2 + 2M\epsilon)$. • *Proof.* We have

$$\tilde{\boldsymbol{\theta}}^{\top}(\hat{\boldsymbol{\Gamma}} - \boldsymbol{\Gamma}^*)\tilde{\boldsymbol{\theta}} = \sum_{t_1=1}^{p^2} \sum_{t_2=1}^{p^2} \tilde{\boldsymbol{\theta}}_{Gt_1}^{\top}(\boldsymbol{\Delta}_{\Gamma})_{Gt_1,Gt_2} \tilde{\boldsymbol{\theta}}_{Gt_2}.$$
 (135)

Therefore,

$$|\tilde{\boldsymbol{\theta}}^{\top}(\hat{\boldsymbol{\Gamma}} - \boldsymbol{\Gamma}^{*})\tilde{\boldsymbol{\theta}}| \leq \sum_{t_{1}=1}^{p^{2}} \sum_{t_{2}=1}^{p^{2}} |\tilde{\boldsymbol{\theta}}_{Gt_{1}}^{\top}(\boldsymbol{\Delta}_{\Gamma})_{Gt_{1},Gt_{2}}\tilde{\boldsymbol{\theta}}_{Gt_{2}}|$$

$$\leq \sum_{t_{1}=1}^{p^{2}} \sum_{t_{2}=1}^{p^{2}} ||\tilde{\boldsymbol{\theta}}_{Gt_{1}}||_{2} ||(\boldsymbol{\Delta}_{\Gamma})_{Gt_{1},Gt_{2}}||_{F} ||\tilde{\boldsymbol{\theta}}_{Gt_{2}}||_{2}$$

$$\leq ||\boldsymbol{\mathcal{C}}(\boldsymbol{\Delta}_{\Gamma})||_{\infty} \sum_{t_{1}=1}^{p^{2}} \sum_{t_{2}=1}^{p^{2}} ||\tilde{\boldsymbol{\theta}}_{Gt_{1}}||_{2} ||\tilde{\boldsymbol{\theta}}_{Gt_{2}}||_{2}$$

$$\leq (\epsilon^{2} + 2M\epsilon) ||\tilde{\boldsymbol{\theta}}||_{\bar{G},2}^{2}, \qquad (136)$$

where we used (117). We have

$$\|\tilde{\boldsymbol{\theta}}\|_{\bar{\mathcal{G}},2}^{2} = \|\tilde{\boldsymbol{\theta}}_{\mathcal{M}} + \tilde{\boldsymbol{\theta}}_{\mathcal{M}^{\perp}}\|_{\bar{\mathcal{G}},2}^{2} = (\|\tilde{\boldsymbol{\theta}}_{\mathcal{M}}\|_{\bar{\mathcal{G}},2} + \|\tilde{\boldsymbol{\theta}}_{\mathcal{M}^{\perp}}\|_{\bar{\mathcal{G}},2})^{2}$$

$$\stackrel{(129)}{\leq} 16 \|\tilde{\boldsymbol{\theta}}_{\mathcal{M}}\|_{\bar{\mathcal{G}},2}^{2} \stackrel{(122)}{\leq} 16s_{n} \|\tilde{\boldsymbol{\theta}}_{\mathcal{M}}\|_{2}^{2} \leq 16s_{n} \|\tilde{\boldsymbol{\theta}}\|_{2}^{2}. \quad (137)$$

Noting that $\tilde{\theta}^{\top} \Gamma^*)\tilde{\theta} \ge \phi^*_{min} \|\tilde{\theta}\|_2^2$ and using (126), (136) and (137), we have

$$\delta \mathcal{L}(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \ge \left(\frac{1}{2} \phi_{min}^* - 8s_n(\epsilon^2 + 2M\epsilon)\right) \|\tilde{\boldsymbol{\theta}}\|_2^2 = \kappa_{\mathcal{L}} \|\tilde{\boldsymbol{\theta}}\|_2^2,$$

proving the desired result. \square

Proof of Theorem 2. First choose ϵ to make $\kappa_{\mathcal{L}} > 0$ in Lemma 6. For instance, suppose we take $8s_n(\epsilon^2 + 2M\epsilon) \le \phi_{min}^*/4$. Then $\kappa_{\mathcal{L}} > \phi_{min}^*/4$. Now pick

$$\epsilon = C_0 \sqrt{ln(p_n)/n} \le \min\left\{M, \frac{\phi_{min}^*}{96s_n M}\right\}, \tag{138}$$

leading to $8s_n(\epsilon^2 + 2M\epsilon) \le 24s_nM\epsilon \le \phi_{min}^*/4$. These upper bounds can be ensured by picking appropriate lower bounds to sample size n and invoking Lemma 2. The choice of n specified in (26) satisfies (138) with probability $> 1 - 2/p_n^{\tau-2}$. Using $\epsilon = C_0 \sqrt{ln(p_n)/n} \le M$, the lower bound on λ_n given

in (25) satisfies (130) with $\mathcal{R}^*(\nabla \mathcal{L}(\theta^*))$ as in Lemma 5. By [30, Theorem 1], $\hat{\theta}$ given by (127) satisfies

$$\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*\|_2 \le \frac{3\lambda_n}{\kappa_{\mathcal{L}}} \Psi(\mathcal{M}). \tag{139}$$

The left-side of (139) equals $\|\hat{\Delta} - \Delta^*\|_F$ while the right-side of (139) equals right-side of (27) using $\Psi(\mathcal{M}) = \sqrt{s_n}$, $\kappa_{\mathcal{L}} > \phi_{min}^*/4$, and noting that $\max_{t=1,\dots,p^2} \|\boldsymbol{\theta}_{Gt}^*\|_2 = \max_{\{k,\ell\}\in V\times V} \|(\boldsymbol{\Delta}^*)^{(k\ell)}\|_F$. This proves Theorem 2.

REFERENCES

- S.L. Lauritzen, Graphical models. Oxford, UK: Oxford Univ. Press, 1996.
- [2] J. Whittaker, Graphical Models in Applied Multivariate Statistics. New York: Wiley, 1990.
- [3] P. Danaher, P. Wang and D.M. Witten, "The joint graphical lasso for inverse covariance estimation across multiple classes," *J. Royal Statistical Society, Series B (Methodological)*, vol. 76, pp. 373-397, 2014.
- [4] J. Friedman, T. Hastie and R. Tibshirani, "Sparse inverse covariance estimation with the graphical lasso," *Biostatistics*, vol. 9, no. 3, pp. 432-441, July 2008.
- [5] N. Meinshausen and P. Bühlmann, "High-dimensional graphs and variable selection with the Lasso," *Ann. Statist.*, vol. 34, no. 3, pp. 1436-1462, 2006.
- [6] K. Mohan, P. London, M. Fazel, D. Witten and S.I. Lee, "Node-based learning of multiple Gaussian graphical models," *J. Machine Learning Research*, vol. 15, pp. 445-488, 2014.
- [7] Y. Wu, T. Li, X. Liu and L.I Chen, "Differential network inference via the fused D-trace loss with cross variables," *Electronic J. Statistics*, vol. 14, pp. 1269-1301, 2020.
- [8] H. Yuan, R. Xi, C. Chen and M. Deng, "Differential network analysis via lasso penalized D-trace loss," *Biometrika*, vol. 104, pp. 755-770, 2017
- [9] Z. Tang, Z. Yu and C. Wang, "A fast iteraive algorithm for highdimensional differential network," *Computational Statistics*, vol. 35, pp. 95-109, 2020.
- [10] B. Zhao, Y.S. Wang and M. Kolar, "Direct estimation of differential functional graphical models," in *Proc. 33rd Conf. Neural Information Processing Systems (NeurIPS 2019)*, Vancouver, Canada, 2019.
- [11] B. Zhao, Y.S. Wang and M. Kolar, "FuDGE: A method to estimate a functional differential graph in a high-dimensional setting," *J. Machine Learning Research*, vol. 23, pp. 1-82, 2022.
- [12] S.D. Zhao, T.T. Cai and H. Li, "Direct estimation of differential networks," *Biometrika*, vol. 101, pp. 253-268, June 2014.
- [13] E. Belilovsky, G. Varoquaux and M.B. Blaschko, "Hypothesis testing for differences in Gaussian graphical models: Applications to brain connectivity," *Advances in Neural Information Processing Systems* (NIPS 2016), vol. 29, Dec. 2016.
- [14] J. Chiquet, G. Rigaill and M. Sundquist, "A multiattribute Gaussian graphical model for inferring multiscale regulatory networks: an application in breast cancer." In: Sanguinetti G., Huynh-Thu V. (eds), Gene Regulatory Networks. Methods in Molecular Biology, vol 1883. Humana Press, New York, NY, 2019.
- [15] M. Kolar, H. Liu and E.P. Xing, "Markov network estimation from multi-attribute data," in *Proc. 30th Intern. Conf. Machine Learning* (ICML), Atlanta, GA, 2013.
- [16] M. Kolar, H. Liu and E.P. Xing, "Graph estimation from multi-attribute data," J. Machine Learning Research, vol. 15, pp. 1713-1750, 2014.
- [17] J.K. Tugnait, "Sparse-group lasso for graph learning from multiattribute data," *IEEE Trans. Signal Process.*, vol. 69, pp. 1771-1786, 2021. (Corrections, vol. 69, p. 4758, 2021.)
- [18] G. Marjanovic and V. Solo, "Vector l_0 sparse conditional independence graphs," in *Proc. IEEE ICASSP 2018*, pp. 2731-2735, 2018.
- [19] Z. Yue, P. Sundaram and V. Solo, "Fast block-sparse estimation for vector networks," in *Proc. IEEE ICASSP* 2020, pp. 5505-5509, 2020.
- [20] P. Sundaram, M. Luessi, M. Bianciardi, S. Stufflebeam, M. Hämäläinen and V. Solo, "Individual resting-state brain networks enabled by massive multivariate conditional mutual information," *IEEE Trans. Med. Imaging*, vol. 39, pp. 1957-1966, 2020.
- [21] Z. Yue and V. Solo, "Comparing vector networks via frequency domain persistent homology," in *Proc. IEEE CDC*, pp. 126-131, Dec. 2021.

- [22] K. Tsai, O. Koyejo, and M. Kolar, "Joint Gaussian graphical model estimation: A survey," WIREs Computational Statistics, vol. 14, no. 6, pp. e1582, Nov/Dec. 2022.
- [23] B. Zhang, H. Li, R.B. Riggins, M. Zhan, J. Xuan, Z. Zhang, E.P. Hoffman, R. Clarke and Y. Wang, "Differential dependency network analysis to identify condition-specific topological changes in biological networks," *Bioinformatics*, vol. 25, no. 4, pp. 526-532, 2009.
- [24] P. Xu and Q. Gu, "Semiparametric differential graph models," in *Proc.* 30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain, 2016.
- [25] S. Liu, J.A. Quinn, M.U. Gutmann, T. Suzuki and M. Sugiyama, "Direct learning of sparse changes in Markov networks by density ratio estimation," *Neural Computation*, 26(6):1169-1197, 2014.
- [26] S. Liu, T. Suzuki, R. Relator, J. Sese, M. Sugiyama, and K. Fukumizu, "Support consistency of direct sparse-change learning in Markov networks" *Annals Statistics*, vol. 45, no. 3, pp. 959-990, 2017.
- [27] S. Liu, K. Fukumizu and T. Suzuki, "Learning sparse structural changes in high-dimensional Markov network: A review on methodologies and theories," *Behaviormetrika*, vol. 44, pp. 265-286, 2017.
- [28] B. Jiang, X. Wang and C. Leng, "A direct approach for sparse quadratic discriminant analysis," *J. Machine Learning Research*, vol. 19, pp. 1-37, 2018
- [29] T. Zhang and H. Zou, "Sparse precision matrix estimation via lasso penalized D-trace loss," *Biometrika*, vol. 101, pp. 103-120, 2014.
- [30] S.N. Negahban, P. Ravikumar, M.J. Wainwright and B. Yu, "A unified framework for high-dimensional analysis of M-estimators with decomposable regularizers," *Statistical Science*, vol. 27, No. 4, pp. 538-557, 2012
- [31] P. Ravikumar, M.J. Wainwright, G. Raskutti and B. Yu, "High-dimensional covariance estimation by minimizing ℓ₁-penalized log-determinant divergence," *Electronic J. Statistics*, vol. 5, pp. 935-980, 2011
- [32] J.K. Tugnait, "Estimation of high-dimensional differential graphs from multi-attribute data," in *Proc. 2023 IEEE Intern. Conf. Acoustics, Speech & Signal Processing (ICASSP 2023)*, pp. 1-5, Rhodes Island, Greece, June 4-9, 2023.
- [33] D.S. Tracy and K.G. Jinadasa, "Partitioned Kronecker products of matrices and applications," *Canadian J. Statistics*, vol. 17, pp. 107-120, March 1989.
- [34] S. Liu, "Matrix results on Khatri-Rao and Tracy-Singh products," Linear Algebra & Its Applications, vol. 289, pp. 267-277, 1999.
- [35] M. Yuan and Y. Lin, "Model selection and estimation in regression with grouped variables," J. Royal Statistical Society, Series B (Methodological), vol. 68, pp. 49-67, 2006.
- [36] J. Friedman, T. Hastie and R. Tibshirani, "A note on the group lasso and a sparse group lasso," arXiv:1001.0736v1 [math.ST], 5 Jan 2010.
- [37] N. Simon, J. Friedman, T. Hastie and R. Tibshirani, "A sparse-group lasso," J. Computational Graphical Statistics, vol. 22, pp. 231-245, 2013.
- [38] S. Boyd, N. Parikh, E. Chu, B. Peleato and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends in Machine Learning*, vol. 3, no. 1, pp. 1-122, 2010.
- [39] A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," SIAM Journal on Imaging Sciences, vol. 2, no. 1, pp. 183-202, 2009.
- [40] C. Leng and C.Y. Tang, "Sparse matrix graphical models," J. American Statistical Association, vol. 107, pp. 1187-1200, Sep. 2012.
- [41] T. Tsiligkaridis, A.O. Hero, III, and S. Zhou, "On convergence of Kronecker graphical lasso algorithms," *IEEE Trans. Signal Process.*, vol. 61, no. 7, pp. 1743-1755, April 2013.
- [42] X. Lyu, W.W. Sun, Z. Wang, H. Liu, J. Yang and G. Cheng, "Tensor graphical model: Non-convex optimization and statistical inference," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 42, no. 8, pp. 2024-2037, 1 Aug. 2020.
- [43] A-L. Barabási and R. Albert, Réka, "Emergence of scaling in random networks," *Science*, vol. 286, no. 5439, pp. 509-512, Oct. 1999.
- [44] S. Lu, J. Kang, W. Gong and D. Towsley, "Complex network comparison using random walks," in WWW '14 Companion: Proc. 23rd Intern. Conf. World Wide Web, pp. 727-730, Seoul, Korea, April 2014.
- [45] S. Zhang, B. Guo, A. Dong, J. He, Z. Xu and S.X. Chen, "Cautionary tales on air-quality improvement in Beijing," *Proc. Royal Society* A: Mathematical, Physical and Engineering Sciences, vol. 473, p. 20170457, 2017.
- [46] W. Chen, F. Wang, G. Xiao, J. Wu and S. Zhang, "Air quality of Beijing and impacts of the new ambient air quality standard," *Atmosphere*, vol. 6, pp. 1243-1258, 2015.