

Learning Sparse High-Dimensional Matrix-Valued Graphical Models From Dependent Data

Jitendra K. Tugnait

Abstract—We consider the problem of inferring the conditional independence graph (CIG) of a sparse, high-dimensional, stationary matrix-variate Gaussian time series. All past work on high-dimensional matrix graphical models assumes that independent and identically distributed (i.i.d.) observations of the matrix-variate are available. Here we allow dependent observations. We consider a sparse-group lasso-based frequency-domain formulation of the problem with a Kronecker-decomposable power spectral density (PSD), and solve it via an alternating direction method of multipliers (ADMM) approach. The problem is bi-convex which is solved via flip-flop optimization. We provide sufficient conditions for local convergence in the Frobenius norm of the inverse PSD estimators to the true value. This result also yields a rate of convergence. We illustrate our approach using numerical examples utilizing both synthetic and real data.

Index Terms—Sparse graph learning; matrix graph estimation; matrix time series; undirected graph; inverse spectral density estimation.

I. INTRODUCTION

IN graphical models, graphs display the conditional independence structure of the variables, and learning the graph structure is equivalent to learning a factorization of the joint probability distribution of these random variables [1]. In a vector graphical model, the conditional statistical dependency structure among p random variables x_1, x_2, \dots, x_p , is represented using an undirected graph $\mathcal{G} = (V, \mathcal{E})$ with a set of p vertices (nodes) $V = \{1, 2, \dots, p\} = [p]$, and a corresponding set of (undirected) edges $\mathcal{E} \subseteq [p] \times [p]$. There is no edge between nodes i and j iff x_i and x_j are conditionally independent given the remaining $p-2$ variables. Suppose $\mathbf{x} \sim \mathcal{N}_r(\mathbf{m}, \Sigma)$, with $\mathbf{m} \in \mathbb{R}^p$, $\Sigma \in \mathbb{R}^{p \times p}$, positive definite $\Sigma = \Omega^{-1}$, where $\mathcal{N}_r(\mathbf{m}, \Sigma)$ denotes a real-valued Gaussian vector with mean \mathbf{m} and covariance Σ . Then Ω_{ij} , the (i, j) -th element of Ω , is zero iff x_i and x_j are conditionally independent [1]. Of much interest is the high-dimensional case where p is greater than or of the order of the data sample size n [2]. In particular, in a high-dimensional setting, as $n \uparrow \infty$, $p/n \rightarrow c > 0$, instead of $p/n \rightarrow 0$ as in classical low-dimensional statistical analysis framework [2, Chapter 1]. Such models for \mathbf{x} have been extensively studied [2]–[5]. In this paper we address the problem of high-dimensional matrix graph estimation. If $p/n \ll 1$, we use the term low-dimensional for such cases in this paper.

Consider a stationary p -dimensional multivariate Gaussian time series $\mathbf{x}(t)$, $t = 0, \pm 1, \pm 2, \dots$, with i th component $x_i(t)$. In the corresponding time series graph $\mathcal{G} = (V, \mathcal{E})$, there is no edge between nodes i and j iff $\{x_i(t)\}$ and $\{x_j(t)\}$ are conditionally independent given the remaining $p-2$ scalar

series $\{x_\ell(t), \ell \in [p], \ell \neq i, \ell \neq j\}$ [6]. Denote the power spectral density (PSD) matrix of zero-mean $\{\mathbf{x}(t)\}$ by $\mathbf{S}_x(f)$, where $\mathbf{S}_x(f) = \sum_{\tau=-\infty}^{\infty} \mathbf{R}_{xx}(\tau) e^{-i2\pi f\tau}$, $\mathbf{R}_{xx}(\tau) = E\{\mathbf{x}(t+\tau)\mathbf{x}^\top(t)\}$ and $\iota = \sqrt{-1}$. In [6] it was shown that conditional independence of two time series components given all other components of the zero-mean time series, is encoded by zeros in the inverse PSD, that is, $\{i, j\} \notin \mathcal{E}$ iff the (i, j) -th element of $\mathbf{S}_x^{-1}(f)$, $[\mathbf{S}_x^{-1}(f)]_{ij} = 0$ for every f . In [6] the low-dimensional case is addressed whereas nonparametric frequency-domain approaches for graph estimation in high-dimensional settings have been considered in [7]–[9]. Refs. [7], [9] provide performance analysis and guarantees. Parametric modeling based approaches in low-dimensional settings for conditional independence graph (CIG) estimation for time series are discussed in [10]–[15]. These papers are focused on algorithm development and they do not provide performance guarantees (such as [9, Theorem 1]). Estimation of sparse high-dimensional parametric time series models is discussed in [16] where performance analysis in high-dimensions is carried out, but the graphical modeling aspect is not addressed.

The need for matrix-valued graphical models arises in several applications [17]–[27] (see also related work of [28]). Here we observe matrix-valued time series $\{\mathbf{Z}(t)\}$ where $\mathbf{Z}(t) \in \mathbb{R}^{p \times q}$. If one vectorizes using $\text{vec}(\mathbf{Z})$ where $\text{vec}(\mathbf{Z}) \in \mathbb{R}^{pq}$ denotes column-wise vectorization of \mathbf{Z} , then use of $\text{vec}(\mathbf{Z})$ will result in a pq -node graph with $(pq) \times (pq)$ precision matrix, which could be ultra-high-dimensional, and it ignores any structural information among rows and columns of $\mathbf{Z}(t)$ [17]. With \otimes denoting the matrix Kronecker product, the basic idea in matrix-valued graphs is to model the covariance of $\text{vec}(\mathbf{Z})$ as $\Psi \otimes \Sigma$ with $\Psi \in \mathbb{R}^{q \times q}$ and $\Sigma \in \mathbb{R}^{p \times p}$, reducing the number of unknowns from $\mathcal{O}(p^2 q^2)$ to $\mathcal{O}(p^2 + q^2)$, while also preserving the structural information. Given data, one estimates two precision matrices $\Omega = \Sigma^{-1}$ and $\Upsilon = \Psi^{-1}$. In the matrix graph, conditional independence between \mathbf{Z}_{ij} and $\mathbf{Z}_{k\ell}$ is determined by zeros in Ω and Υ [17]. This is the Kronecker graph model [29], [30]: If \mathcal{G}_1 and \mathcal{G}_2 are graphs with adjacency matrices $\mathcal{A}(\mathcal{G}_1)$ and $\mathcal{A}(\mathcal{G}_2)$, respectively, then the Kronecker product graph (KPG) $\mathcal{G}_1 \otimes \mathcal{G}_2$ is defined as the graph with adjacency matrix $\mathcal{A}(\mathcal{G}_1) \otimes \mathcal{A}(\mathcal{G}_2)$ [30, Def. 1]. In our context the nonzero entries of Υ and Ω determine the nonzero entries of the adjacency matrices of graphs \mathcal{G}_1 and \mathcal{G}_2 , respectively, with KPG $\mathcal{G} = \mathcal{G}_1 \otimes \mathcal{G}_2$.

Our objective in this paper is to learn a conditional independence KPG associated with time-dependent matrix-valued zero-mean $p \times q$ Gaussian sequence $\mathbf{Z}(t)$, under high-dimensional settings, given observations of $\{\mathbf{Z}(t)\}_{t=0}^{n-1}$.

A. Related Work

Prior work on KPG estimation under high-dimensional settings [17]–[27] all assume that i.i.d. observations of \mathbf{Z}

J.K. Tugnait is with the Dept. of Elec. & Comp. Eng., 200 Broun Hall, Auburn University, Auburn, AL 36849, USA. Email: tugnajk@auburn.edu . Supported by NSF Grants ECCS-2040536 and CCF-2308473.

are available for graphical modeling. Refs. [17], [18], [25] all solve the same bi-convex optimization problem, using an identical alternating minimization approach, but they differ in theoretical analysis. Ref. [21] uses a Kronecker sum model whereas we use a Kronecker-product separable covariance structure (see (4) later) for $\{\mathbf{Z}(t)\}$.

There is no prior reported work on high-dimensional matrix graph estimation with dependent data using a nonparametric approach. Parametric models using state-space models are estimated in [31], [32] for KPG estimation in low-dimensional settings using Bayesian approaches. Granger causality graphs (not the same as CIGs) for matrix time series are estimated in [33] using first-order AR models and in [34] using an infinite dimensional model class (which includes ARMAX models of any order), both in low-dimensional settings (i.e., $pq/n \ll 1$ and/or $\lim_{n \rightarrow \infty} pq/n = 0$, with pq representing number of nodes in KPG). In contrast, this paper considers conditional independence KPG's under high-dimensional settings. Ref. [35] investigates sum of Kronecker product AR models for matrix times series with no consideration of CIGs. Estimation of a KPG model corresponding to an AR Gaussian process is investigated in [36] in low-dimensional settings with no performance analysis or guarantees. A distinguishing aspect of [36] is that it imposes a Kronecker product decomposition on the support of the inverse PSD, not the inverse PSD of the time series. With regard to [34], [36], we note that in the synthetic data example using an ARMA model in [34, Sec. 6.1, Fig. 2], number of nodes is 16 ($= pq$) and sample size is $n = 3900$, leading to $pq/n = 0.004$, a low-dimensional setting. The real data example of [34, Sec. 6.2] does have $pq = 96$ and $n = 500$, implying $pq/n = 0.19$. The distinction is that the ground truth is known in synthetic data examples permitting evaluation of the efficacy of the considered approach, whereas such is not the case in real data examples. Thus [34] does not address the high-dimensional scenario as relatively high $pq/n = 0.19$ in their real data example is not supported by any commensurate synthetic data example. In contrast, we provide such support, as seen in Table I, Sec. VI-A of this paper, where $pq = 225$ and varying $n \in \{64, 128, 256, 512, 1024, 2148\}$, implying $pq/n \in \{3.5, 1.76, 0.88, 0.44, 0.22, 0.11\}$. In our real data example (Sec. VI-B), we have $pq = 88$ and $n = 364$ with $pq/n = 0.24$. In [36, Sec. 6.1], the synthetic data example has $pq = 36$ and $n = 1000$ or 2000 ($pq/n = 0.036$ or 0.018), again a low-dimensional scenario. The real data example of [36, Sec. 6.2] has $pq = 36$ and $n = 389$ ($pq/n = 0.09$). The comments made pertaining to [34] regarding differences in pq/n ratios for real and synthetic data examples, apply to [36] as well. Finally, [37] considers a first-order matrix AR model for matrix time series where when vectorized, the vectorized time-series AR coefficient is expressed as a Kronecker product. Low-dimensional asymptotics are provided in [37] and the issue of the underlying CIG is not addressed.

A frequency-domain formulation is used in this paper, following the approach of [9] for dependent vector time series. The resulting optimization problem is bi-convex, as in [17], [18], [25], but with complex variables, and is solved via an alternating minimization approach using Wirtinger calculus [38] for optimization of real functions of complex variables.

A preliminary version of parts of this paper appear in a workshop paper [39]. Theorems 1-3 and their proofs, and the real data example do not appear in [39].

B. Our Contributions, Outline and Notation

The underlying system model including a generative model (5) for time-dependent matrix Gaussian sequence, is presented in Sec. II. A frequency-domain based penalized log-likelihood objective function is derived in Sec. III for estimation of the matrix graph, resulting in a Kronecker-decomposable power spectral density representation (15). A flip-flop algorithm based on two ADMM algorithms is presented in Sec. IV to optimize the bi-convex objective function. In Sec. V the performance of the proposed optimization algorithm is analyzed under a high-dimensional large sample setting in Theorems 1-3, patterned after [22] and exploiting some results from [9], [40], [41]. Numerical results are presented in Sec. VI and proofs of Theorems 1, 2 and 3 are given in three appendices.

Notation. The superscripts $*$, \top and H denote the complex conjugate, transpose and conjugate transpose operations, respectively, \mathbb{R} and \mathbb{C} denote the sets of real and complex numbers, respectively, and $\text{Re}(\mathbf{x})$ is the real part of $\mathbf{x} \in \mathbb{C}^p$. We use $\iota := \sqrt{-1}$. A $p \times p$ identity matrix is denoted by \mathbf{I}_p . Given $\mathbf{A} \in \mathbb{C}^{p \times p}$, $\phi_{\min}(\mathbf{A})$, $\phi_{\max}(\mathbf{A})$, $|\mathbf{A}|$, $\text{tr}(\mathbf{A})$ and $\text{etr}(\mathbf{A})$ denote the minimum eigenvalue, maximum eigenvalue, determinant, trace, and exponential of trace of \mathbf{A} , respectively. We use $\mathbf{A} \succeq 0$ and $\mathbf{A} \succ 0$ to denote that Hermitian \mathbf{A} is positive semi-definite and positive definite, respectively. For $\mathbf{B} \in \mathbb{C}^{p \times q}$, we define the operator norm, the Frobenius norm and the vectorized ℓ_1 norm, respectively, as $\|\mathbf{B}\| = \sqrt{\phi_{\max}(\mathbf{B}^H \mathbf{B})}$, $\|\mathbf{B}\|_F = \sqrt{\text{tr}(\mathbf{B}^H \mathbf{B})}$ and $\|\mathbf{B}\|_1 = \sum_{i,j} |B_{ij}|$, where B_{ij} is the (i, j) -th element of \mathbf{B} , also denoted by $[\mathbf{B}]_{ij}$. For vector $\boldsymbol{\theta} \in \mathbb{C}^p$, we define $\|\boldsymbol{\theta}\|_1 = \sum_{i=1}^p |\theta_i|$ and $\|\boldsymbol{\theta}\|_2 = \sqrt{\sum_{i=1}^p |\theta_i|^2}$, and we also use $\|\boldsymbol{\theta}\|$ for $\|\boldsymbol{\theta}\|_2$. Given $\mathbf{A} \in \mathbb{C}^{p \times p}$, $\mathbf{A}^+ = \text{diag}(\mathbf{A})$ is a diagonal matrix with the same diagonal as \mathbf{A} , and $\mathbf{A}^- = \mathbf{A} - \mathbf{A}^+$ is \mathbf{A} with all its diagonal elements set to zero. We use \mathbf{A}^{-*} for $(\mathbf{A}^*)^{-1}$, the inverse of complex conjugate of \mathbf{A} , and $\mathbf{A}^{-\top}$ for $(\mathbf{A}^\top)^{-1}$. Given $\mathbf{A} \in \mathbb{C}^{n \times p}$, column vector $\text{vec}(\mathbf{A}) \in \mathbb{C}^{np}$ denotes the vectorization of \mathbf{A} which stacks the columns of the matrix \mathbf{A} . The notation $\mathbf{y}_n = \mathcal{O}_P(\mathbf{x}_n)$ for random $\mathbf{y}_n, \mathbf{x}_n \in \mathbb{C}^p$ means that for any $\varepsilon > 0$, there exists $0 < T < \infty$ such that $P(\|\mathbf{y}_n\| \leq T\|\mathbf{x}_n\|) \geq 1 - \varepsilon \forall n \geq 1$. The notation $\mathbf{x} \sim \mathcal{N}_c(\mathbf{m}, \boldsymbol{\Sigma})$ denotes a complex random vector \mathbf{x} that is circularly symmetric (proper), complex Gaussian with mean \mathbf{m} and covariance $\boldsymbol{\Sigma}$, and $\mathbf{x} \sim \mathcal{N}_r(\mathbf{m}, \boldsymbol{\Sigma})$ denotes real-valued Gaussian \mathbf{x} with mean \mathbf{m} and covariance $\boldsymbol{\Sigma}$.

II. SYSTEM MODEL

A random matrix $\mathbf{Z} \in \mathbb{R}^{p \times q}$ is said to have a matrix normal (Gaussian) distribution if its pdf $f(\mathbf{Z}|\mathbf{M}, \boldsymbol{\Sigma}, \boldsymbol{\Psi})$, characterized by $\mathbf{M} \in \mathbb{R}^{p \times q}$, $\boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$, $\boldsymbol{\Psi} \in \mathbb{R}^{q \times q}$, is [42, Chap. 2]

$$f(\mathbf{Z}|\mathbf{M}, \boldsymbol{\Sigma}, \boldsymbol{\Psi}) = \frac{\text{etr}\left(-\frac{1}{2}(\mathbf{Z} - \mathbf{M})\boldsymbol{\Psi}^{-1}(\mathbf{Z} - \mathbf{M})^\top \boldsymbol{\Sigma}^{-1}\right)}{(2\pi)^{pq/2} |\boldsymbol{\Sigma}|^{q/2} |\boldsymbol{\Psi}|^{p/2}}. \quad (1)$$

We will use the notation $\mathbf{Z} \sim \mathcal{MN}(\mathbf{M}, \boldsymbol{\Sigma}, \boldsymbol{\Psi})$ for the matrix normal distribution specified by (1). Equivalently,

$$\text{vec}(\mathbf{Z}) \sim \mathcal{N}_r(\text{vec}(\mathbf{M}), \boldsymbol{\Psi} \otimes \boldsymbol{\Sigma}). \quad (2)$$

Here Ψ is the row covariance matrix and Σ is the column covariance matrix [42] since the k th column $\mathbf{Z}_{\cdot k} \sim \mathcal{N}_r(\mathbf{0}, [\Psi]_{kk}\Sigma)$ and the i th row $\mathbf{Z}_i^\top \sim \mathcal{N}_r(\mathbf{0}, [\Sigma]_{ii}\Psi)$.

With $\mathbf{Z} \in \mathbb{R}^{p \times q}$ modeled as a zero-mean matrix normal vector and $\mathbf{z} = \text{vec}(\mathbf{Z})$, [17] assumes

$$E\{\mathbf{z}\mathbf{z}^\top\} = \Psi \otimes \Sigma, \quad (3)$$

implying a separable covariance structure [28]. Let $\Omega = \Sigma^{-1}$ and $\Upsilon = \Psi^{-1}$ denote the respective precision matrices. Then \mathbf{Z}_{ij} and $\mathbf{Z}_{k\ell}$ are conditionally independent given remaining entries in \mathbf{Z} iff (i) at least one of Ω_{ik} and $\Upsilon_{j\ell}$ is zero when $i \neq k, j \neq \ell$, (ii) $\Omega_{ik} = 0$ when $i \neq k, j = \ell$, and (iii) $\Upsilon_{j\ell} = 0$ when $i = k, j \neq \ell$ [17].

In this paper we will model our time-dependent zero-mean matrix-valued, stationary, $p \times q$ Gaussian sequence $\mathbf{Z}(t)$, $\mathbf{z}(t) = \text{vec}(\mathbf{Z}(t))$, as having the separable covariance structure given by

$$E\{\mathbf{z}(t+\tau)\mathbf{z}^\top(t)\} = \Psi(\tau) \otimes \Sigma \quad (4)$$

where $\Psi(\tau)$, $\tau = 0, \pm 1, \dots$ models time-dependence while $\Sigma \succ \mathbf{0}$ is fixed. Under (4), the row covariance sequence is $E\{\mathbf{Z}_i^\top(t+\tau)\mathbf{Z}_i(t)\} = [\Sigma]_{ii}\Psi(\tau)$ and the column covariance sequence is $E\{\mathbf{Z}_{\cdot k}(t+\tau)\mathbf{Z}_{\cdot k}^\top(t)\} = \Sigma[\Psi(\tau)]_{kk}$. Thus we allow possible temporal dependence in matrix observations via $\Psi(\tau)$. With $\{\mathbf{e}(t)\}$ i.i.d., $\mathbf{e}(t) \sim \mathcal{N}_r(\mathbf{0}, \mathbf{I}_{pq})$, a generative model for $\mathbf{z}(t)$ is given by

$$\mathbf{z}(t) = \sum_{i=0}^L (\mathbf{B}_i \otimes \mathbf{F})\mathbf{e}(t-i), \quad \mathbf{B}_i \in \mathbb{R}^{q \times q}, \quad \mathbf{F} \in \mathbb{R}^{p \times p} \quad (5)$$

$$\Rightarrow E\{\mathbf{z}(t+\tau)\mathbf{z}^\top(t)\} = \underbrace{\left(\sum_{i=0}^L \mathbf{B}_i \mathbf{B}_i^\top\right)}_{=\Psi(\tau)} \otimes \underbrace{(\mathbf{F}\mathbf{F}^\top)}_{=\Sigma}. \quad (6)$$

In (5), we can have $L \uparrow \infty$ so long as assumption (A2) stated in Sec. III holds. In sequel, we exploit (4) in our approach without considering (6), the latter is used only for synthetic data generation.

The PSD of $\{\mathbf{z}(t)\}$ is $\mathbf{S}_z(f) = \bar{\mathbf{S}}(f) \otimes \Sigma$ where $\bar{\mathbf{S}}(f) = \sum_{\tau} \Psi(\tau)e^{-i2\pi f\tau}$. Then $\mathbf{S}_z^{-1}(f) = \bar{\mathbf{S}}^{-1}(f) \otimes \Sigma^{-1}$, and by [6], in the pq -node graph $\mathcal{G} = (V, \mathcal{E})$, $|V| = pq$, associated with $\{\mathbf{z}(t)\}$, edge $\{i, j\} \notin \mathcal{E}$ iff $[\mathbf{S}_z^{-1}(f)]_{ij} = 0$ for every f . This does not account for the separable structure of our model. Noting that $\bar{\mathbf{S}}^{-1}(f)$, $f \in [0, 0.5]$, plays the role of $\Upsilon = \Psi^{-1}$, using [6], [17] (also [30, Observation 1]), we deduce that $\{\mathbf{Z}_{ij}(t)\}$ and $\{\mathbf{Z}_{k\ell}(t)\}$ are conditionally independent given remaining entries in $\{\mathbf{Z}(t)\}$ iff (i) at least one of Ω_{ik} and $[\bar{\mathbf{S}}^{-1}(f)]_{j\ell}$, for every $f \in [0, 0.5]$, is zero when $i \neq k, j \neq \ell$, (ii) $\Omega_{ik} = 0$ when $i \neq k, j = \ell$, and (iii) $[\bar{\mathbf{S}}^{-1}(f)]_{j\ell} = 0$, for every $f \in [0, 0.5]$ when $i = k, j \neq \ell$. That is, we have a KPG $\mathcal{G} = \mathcal{G}_1 \otimes \mathcal{G}_2$ where the adjacency matrix of \mathcal{G}_1 is specified by the nonzero entries of $\bar{\mathbf{S}}^{-1}(f) \forall f \in [0, 0.5]$, and that of \mathcal{G}_2 follows from the nonzero entries of Ω .

Our objective is to learn the graph associated with $\{\mathbf{Z}(t)\}$ under some sparsity constraints on Ω and $\bar{\mathbf{S}}^{-1}(f)$, $f \in [0, 0.5]$. Since $\alpha\bar{\mathbf{S}}^{-1}(f) \otimes (\alpha^{-1}\Omega) = \bar{\mathbf{S}}^{-1}(f) \otimes \Omega$, to resolve scaling ambiguity, we could normalize $\|\Omega\|_F = 1$ or $\|\bar{\mathbf{S}}^{-1}(f_1) \cdots \bar{\mathbf{S}}^{-1}(f_M)\|_F = 1$ for suitably placed M frequencies in $(0, 0.5)$; we will follow the latter as stated later in step 2 of Sec. IV-A.

III. PENALIZED NEGATIVE LOG-LIKELIHOOD

Given $\mathbf{z}(t)$ for $t = 0, 1, 2, \dots, n-1$. Define the (normalized) DFT's $\mathbf{d}_z(f_m)$ and $\mathbf{D}_z(f_m)$ of $\mathbf{z}(t)$ and $\mathbf{Z}(t)$, respectively, as (recall $\iota = \sqrt{-1}$),

$$\mathbf{d}_z(f_m) = \frac{1}{\sqrt{n}} \sum_{t=0}^{n-1} \mathbf{z}(t) \exp(-i2\pi f_m t), \quad (7)$$

$$\mathbf{D}_z(f_m) = \frac{1}{\sqrt{n}} \sum_{t=0}^{n-1} \mathbf{Z}(t) \exp(-i2\pi f_m t), \quad (8)$$

$$f_m = m/n, \quad m = 0, 1, \dots, n-1. \quad (9)$$

Then $\mathbf{d}_z(f_m) = \text{vec}(\mathbf{D}_z(f_m))$. It is established in [43] (see also [9]) that, for even n , the set of random vectors $\{\mathbf{d}_z(f_m)\}_{m=0}^{n/2}$ is a sufficient statistic for any inference problem based on dataset $\{\mathbf{z}(t)\}_{t=0}^{n-1}$. Suppose $\mathbf{S}_z(f_k)$ is locally smooth, so that $\mathbf{S}_z(f_k)$ is (approximately) constant over $K = 2m_t + 1$ consecutive frequency points f_m 's where m_t is the half-window size; in our case, this assumption applies to $\bar{\mathbf{S}}(f_k)$. Pick $M = \lfloor (\frac{n}{2} - m_t - 1)/K \rfloor$ and

$$\tilde{f}_k = ((k-1)K + m_t + 1)/n, \quad k \in [M], \quad (10)$$

yielding M equally spaced frequencies \tilde{f}_k in the interval $(0, 0.5)$. We state the local smoothness assumption as assumption (A1).

(A1) Assume that for $\ell = -m_t, -m_t + 1, \dots, m_t$,

$$\mathbf{S}_z(\tilde{f}_{k,\ell}) = \mathbf{S}_z(\tilde{f}_k), \quad (11)$$

$$\text{where } \tilde{f}_{k,\ell} = ((k-1)K + m_t + 1 + \ell)/n. \quad (12)$$

We will invoke [44, Theorem 4.4.1] for distribution of $\mathbf{d}_z(f_m)$. To this end we need assumption (A2).

(A2) The matrix time series $\{\mathbf{Z}(t)\}_{t=-\infty}^{\infty}$ is zero-mean stationary, Gaussian, satisfying $\sum_{\tau=-\infty}^{\infty} |[\Psi(\tau)]_{k\ell}| < \infty$ for every $k, \ell \in [q]$.

By [44, Theorem 4.4.1], under assumption (A1), asymptotically (as $n \rightarrow \infty$), $\mathbf{d}_z(f_m)$, $m = 1, 2, \dots, (n/2)-1$, (n even), are independent proper, complex Gaussian $\mathcal{N}_c(\mathbf{0}, \mathbf{S}_z(f_m))$ random vectors, respectively. Denote the joint probability density function of $\mathbf{d}_z(f_m)$, $m = 1, 2, \dots, (n/2)-1$, as $f_{\mathcal{D}}(\mathcal{D})$ where $\mathcal{D} = \{\mathbf{D}_z(f_m)\}_{m=1}^{(n/2)-1}$. Then we have [9], [43]

$$f_{\mathcal{D}}(\mathcal{D}) = \prod_{k=1}^M \left[\prod_{\ell=-m_t}^{m_t} \frac{\exp(-g_{kl} - g_{kl}^*)}{\pi^{pq} |\mathbf{B}_k|^{1/2} |\mathbf{B}_k^*|^{1/2}} \right], \quad (13)$$

$$g_{kl} = \frac{1}{2} \mathbf{d}_z^H(\tilde{f}_{k,\ell}) (\bar{\mathbf{S}}^{-1}(\tilde{f}_k) \otimes \Sigma^{-1}) \mathbf{d}_z(\tilde{f}_{k,\ell}), \quad (14)$$

$$\mathbf{B}_k = \bar{\mathbf{S}}(\tilde{f}_k) \otimes \Sigma. \quad (15)$$

Using $\text{tr}(\mathbf{A}^\top \mathbf{B} \mathbf{C} \mathbf{G}^\top) = (\text{vec}(\mathbf{A}))^\top (\mathbf{G} \otimes \mathbf{B}) \text{vec}(\mathbf{C})$ and parametrizing in terms of $\Phi_k := \bar{\mathbf{S}}^{-1}(\tilde{f}_k)$ and $\Omega = \Sigma^{-1}$, we have

$$\begin{aligned} g_{kl} &= \frac{1}{2} \mathbf{D}_z^H(\tilde{f}_{k,\ell}) \Sigma^{-1} \mathbf{D}_z(\tilde{f}_{k,\ell}) (\bar{\mathbf{S}}^{-1}(\tilde{f}_k))^\top \\ &= \frac{1}{2} \mathbf{D}_z^H(\tilde{f}_{k,\ell}) \Omega \mathbf{D}_z(\tilde{f}_{k,\ell}) \Phi_k^*. \end{aligned} \quad (16)$$

Define $q \times (qM)$ matrix Γ as

$$\Gamma = [\Phi_1 \ \Phi_2 \ \cdots \ \Phi_M]. \quad (17)$$

Using $|\mathbf{B}_k| = |\bar{\mathbf{S}}(\tilde{f}_k) \otimes \Sigma| = |\bar{\mathbf{S}}(\tilde{f}_k)|^p |\Sigma|^q$, up to some constants the negative log-likelihood follows from (13) as

$$\begin{aligned} -\frac{1}{KMpq} \ln f_{\mathcal{D}}(\mathcal{D}) &\propto G(\Omega, \Gamma, \Gamma^*) \\ &:= -\frac{1}{p} \ln(|\Omega|) - \frac{1}{2Mq} \sum_{k=1}^M (\ln(|\Phi_k|) + \ln(|\Phi_k^*|)) \\ &\quad + \frac{1}{2Mq} \sum_{k=1}^M \text{tr}(\mathbf{A}_k + \mathbf{A}_k^*), \end{aligned} \quad (18)$$

$$\mathbf{A}_k = \frac{1}{Kp} \sum_{\ell=-m_t}^{m_t} \mathbf{D}_z^H(\tilde{f}_{k,\ell}) \Omega \mathbf{D}_z(\tilde{f}_{k,\ell}) \Phi_k^*. \quad (19)$$

In the high-dimension case, to enforce sparsity and to make the problem well-conditioned, we propose to minimize a penalized version $\mathcal{L}(\Omega, \Gamma)$ w.r.t. Ω and Γ ,

$$\mathcal{L}(\Omega, \Gamma) = G(\Omega, \Gamma, \Gamma^*) + P_p(\Omega) + P_q(\{\Phi\}), \quad (20)$$

$$P_p(\Omega) = \lambda_p \sum_{i \neq j}^p |\Omega_{ij}| = \lambda_p \|\Omega^-\|_1, \quad (21)$$

$$\begin{aligned} P_q(\{\Phi\}) &= \alpha \lambda_q \sum_{k=1}^M \sum_{i \neq j}^q [|\Phi_k]_{ij}| \\ &\quad + (1 - \alpha) \sqrt{M} \lambda_q \sum_{i \neq j}^q \|\Phi^{(ij)}\|, \end{aligned} \quad (22)$$

$$\Phi^{(ij)} := [[\Phi_1]_{ij} \ [\Phi_2]_{ij} \ \cdots \ [\Phi_M]_{ij}]^\top \in \mathbb{C}^M, \quad (23)$$

where $\{\Phi\} := \{\Phi_k\}_{k=1}^M$, $\alpha \in [0, 1]$, $\lambda_p, \lambda_q > 0$ are tuning parameters, $P_p(\Omega)$ is the lasso constraint, $P_q(\{\Phi\})$ is a sparse-group lasso sparsity constraint (cf. [45]–[47]) and \sqrt{M} in $P_q(\{\Phi\})$ reflects number of group variables [47].

IV. OPTIMIZATION

The objective function $\mathcal{L}(\Omega, \Gamma)$ in (20) is biconvex: (strictly) convex in Γ , $\Phi_k \succ \mathbf{0}$, for fixed Ω , and (strictly) convex in Ω , $\Omega \succ \mathbf{0}$, for fixed Γ . As is a general approach for biconvex function optimization [48], we will use an iterative and alternating minimization approach where we optimize w.r.t. Ω with Γ fixed, and then optimize w.r.t. Γ with Ω fixed at the last optimized value, and repeat the two optimizations (flip-flop). The algorithm is only guaranteed to converge to a local stationary point of $\mathcal{L}(\Omega, \Gamma)$ [48, Sec. 4.2.1].

With $\hat{\Gamma} = [\hat{\Phi}_1 \ \hat{\Phi}_2 \ \cdots \ \hat{\Phi}_M]$ denoting the estimate of Γ , fix $\Gamma = \hat{\Gamma}$ and let $\mathcal{L}_1(\Omega)$ denote $\mathcal{L}(\Omega, \hat{\Gamma})$ up to some irrelevant constants. We minimize $\mathcal{L}_1(\Omega)$ w.r.t. Ω to estimate $\hat{\Omega}$, where

$$\mathcal{L}_1(\Omega) = -\frac{1}{p} \ln(|\Omega|) + \frac{1}{p} \text{tr}(\Omega \check{\Theta}) + P_p(\Omega), \quad (24)$$

$$\check{\Theta} = \frac{1}{MKq} \sum_{k=1}^M \sum_{\ell=-m_t}^{m_t} \text{Re}\{\mathbf{D}_z(\tilde{f}_{k,\ell}) \hat{\Phi}_k^* \mathbf{D}_z^H(\tilde{f}_{k,\ell})\}. \quad (25)$$

Fix $\Omega = \hat{\Omega}$ and let $\mathcal{L}_2(\Gamma)$ denote $\mathcal{L}(\hat{\Omega}, \Gamma)$ up to some irrelevant constants. We minimize $\mathcal{L}_2(\Gamma)$ w.r.t. Γ to obtain estimate $\hat{\Gamma}$, where

$$\begin{aligned} \mathcal{L}_2(\Gamma) &= -\frac{1}{2Mq} \sum_{k=1}^M (\ln(|\Phi_k|) + \ln(|\Phi_k^*|)) \\ &\quad + \frac{1}{2Mq} \sum_{k=1}^M \text{tr}(\tilde{\Theta}_k \Phi_k + \tilde{\Theta}_k^* \Phi_k^*) + P_q(\{\Phi\}), \end{aligned} \quad (26)$$

$$\tilde{\Theta}_k = \frac{1}{Kp} \sum_{\ell=-m_t}^{m_t} \mathbf{D}_z^\top(\tilde{f}_{k,\ell}) \hat{\Omega} \mathbf{D}_z^*(\tilde{f}_{k,\ell}). \quad (27)$$

Our optimization algorithm is as in Sec. IV-A.

A. Flip-Flop Optimization

1. Initialize $m = 1$, $\Omega^{(0)} = \mathbf{I}_p$, $\Phi_k^{(0)} = \mathbf{I}_q$, $k \in [M]$.
2. Set $\hat{\Omega} = \Omega^{(m-1)}$ in (27). Use the iterative ADMM algorithm [49], as outlined in [9, Sec. 4] and based on Wirtinger calculus [38], to minimize $\mathcal{L}_2(\Gamma)$ (given by (26)) w.r.t. Γ to obtain estimates $\Phi_k^{(m)}$, $k \in [M]$, the M component matrices of the estimate $\Gamma^{(m)}$. Details are in Sec. IV-B and step II of Sec. IV-D. Normalize $\Gamma^{(m)} \leftarrow \Gamma^{(m)} / \|\Gamma^{(m)}\|_F$ to resolve the scaling ambiguity. Let $m \leftarrow m + 1$.
3. Set $\hat{\Gamma} = \Gamma^{(m)}$ in (25). Use the ADMM algorithm of [40, Sec. III] (with $\alpha = 1$ therein, no group-lasso penalty) to minimize $\mathcal{L}_1(\Omega)$ w.r.t. Ω , to obtain estimate $\Omega^{(m)}$. Details are in Sec. IV-C and step IV of Sec. IV-D.
4. Repeat steps 2 and 3 until convergence.

B. ADMM for Estimation of Γ

After variable splitting, the scaled augmented Lagrangian for minimization of $\mathcal{L}_2(\Gamma)$ is [9]

$$\begin{aligned} \mathcal{L}_2^{aL}(\{\Phi\}, \{\mathbf{W}\}, \{\mathbf{U}\}) &= \frac{1}{2Mq} \sum_{k=1}^M \text{tr}(\tilde{\Theta}_k \Phi_k + \tilde{\Theta}_k^* \Phi_k^*) \\ &\quad - \frac{1}{2Mq} \sum_{k=1}^M (\ln(|\Phi_k|) + \ln(|\Phi_k^*|)) + P_q(\{\mathbf{W}\}) \\ &\quad + \frac{\rho}{2} \sum_{k=1}^M \|\Phi_k - \mathbf{W}_k + \mathbf{U}_k\|_F^2 \end{aligned}$$

where $\{\mathbf{U}\} = \{\mathbf{U}_k\}_{k=1}^M$ are dual variables, similarly $\{\mathbf{W}_k\}_{k=1}^M$ are the ‘‘split’’ variables, $\rho > 0$ is the penalty parameter, $\mathbf{U}_k, \mathbf{W}_k \in \mathbb{C}^{q \times q}$. Given the results $\{\tilde{\Phi}^{(i)}\}, \{\mathbf{W}^{(i)}\}, \{\mathbf{U}^{(i)}\}$ of the i th iteration, in the $(i + 1)$ st iteration, the ADMM algorithm executes the following three updates, given in Steps (a)–(c), until convergence. To distinguish between the estimates $\Gamma^{(m)}$ and $\Phi_k^{(m)}$ of the m th iteration of the flip-flip optimization and the estimate of the i th iteration of the ADMM algorithm, we use $\tilde{\Phi}_k^{(i)}$ for the latter.

Step (a). $\{\tilde{\Phi}^{(i+1)}\} \leftarrow \arg \min_{\{\Phi\}} \mathcal{L}_2^{aL}(\{\Phi\}, \{\mathbf{W}^{(i)}\}, \{\mathbf{U}^{(i)}\})$. Up to some terms not dependent upon $\tilde{\Phi}_k$'s [9]

$$\begin{aligned} \mathcal{L}_2^{aL}(\{\Phi\}, \{\mathbf{W}^{(i)}\}, \{\mathbf{U}^{(i)}\}) \\ = \frac{1}{2Mq} \sum_{k=1}^M \mathcal{L}_{2k}^{aL}(\Phi_k, \mathbf{W}_k^{(i)}, \mathbf{U}_k^{(i)}), \end{aligned}$$

$$\begin{aligned} \mathcal{L}_{2k}^{aL}(\Phi_k, \mathbf{W}_k^{(i)}, \mathbf{U}_k^{(i)}) &= -\ln(|\Phi_k|) - \ln(|\Phi_k^*|) \\ &+ \text{tr}(\tilde{\Theta}_k \Phi_k + \tilde{\Theta}_k^* \Phi_k^*) + Mq\rho \|\Phi_k - \mathbf{W}_k^{(i)} + \mathbf{U}_k^{(i)}\|_F^2, \end{aligned}$$

that is, the objective function is separable in k . For each k , the solution is as follows [9]. Let $\mathbf{P}\Delta\mathbf{P}^H$ denote the eigen-decomposition of the Hermitian $\tilde{\Theta}_k - Mq\rho(\mathbf{W}_k^{(i)} - \mathbf{U}_k^{(i)})$ with diagonal matrix Δ consisting of the eigenvalues and $\mathbf{P}\mathbf{P}^H = \mathbf{P}^H\mathbf{P} = \mathbf{I}_q$. Then $\tilde{\Phi}_k^{(i+1)} = \mathbf{P}\tilde{\Delta}\mathbf{P}^H$ where $\tilde{\Delta}$ is the diagonal matrix with ℓ th diagonal element

$$[\tilde{\Delta}]_{\ell\ell} = \frac{-[\Delta]_{\ell\ell} + \sqrt{([\Delta]_{\ell\ell})^2 + 4Mq\rho}}{2Mq\rho}.$$

Step (b). Here we have

$$\{\mathbf{W}^{(i+1)}\} \leftarrow \arg \min_{\{\mathbf{W}\}} \mathcal{L}_2^{aL}(\{\tilde{\Phi}^{(i+1)}\}, \{\mathbf{W}\}, \{\mathbf{U}^{(i)}\}).$$

We update $\{\mathbf{W}_k^{(i+1)}\}_{k=1}^M$ as the minimizer w.r.t. $\{\mathbf{W}\}_{k=1}^M$ of

$$\frac{\rho}{2} \sum_{k=1}^M \|\mathbf{W}_k - (\tilde{\Phi}_k^{(i+1)} + \mathbf{U}_k^{(i)})\|_F^2 + P_q(\{\mathbf{W}\}).$$

The solution follows from [9, Lemma 1]. Let $\mathbf{G}_k = \tilde{\Phi}_k^{(i+1)} + \mathbf{U}_k^{(i)} \in \mathbb{C}^{q \times q}$ and let $\mathbf{G}^{(j\ell)} \in \mathbb{C}^M$ be defined as in (23), but based on \mathbf{G}_k 's. Then the update of $\{\mathbf{W}\}$ is given by

$$\begin{aligned} [\mathbf{W}_k^{(i+1)}]_{j\ell} &= [\mathbf{G}_k]_{jj}, \quad \text{if } j = \ell \\ [\mathbf{W}_k^{(i+1)}]_{j\ell} &= \left(1 - \frac{(1-\alpha)\lambda_q\sqrt{M}}{\rho\|\mathbf{S}_F(\mathbf{G}^{(j\ell)}), \alpha\lambda_q/\rho\|} \right) + \\ &\times S_F\left([\mathbf{G}^{(j\ell)}]_k, \frac{\alpha\lambda_q}{\rho}\right), \quad \text{if } j \neq \ell, \end{aligned}$$

where $(b)_+ := \max(0, b)$, $S_F(b, \beta) := (1 - \beta/|b|)_+ b$ (for complex scalar $b \neq 0$) is the soft-thresholding scalar operator, and $[\mathbf{S}_F(\mathbf{a}, \beta)]_j = S(a_j, \beta)$ with $a_j = [\mathbf{a}]_j$, is the soft-thresholding vector operator.

Step (c). $\{\mathbf{U}^{(i+1)}\} \leftarrow \{\mathbf{U}^{(i)}\} + (\{\tilde{\Phi}^{(i+1)}\} - \{\mathbf{W}^{(i+1)}\})$.

C. ADMM for Estimation of Ω

Using variable splitting, consider

$$\min_{\Omega > 0, \bar{\mathbf{W}}} \left\{ \frac{1}{p} (\text{tr}(\tilde{\Theta}\Omega) - \ln(|\Omega|)) + \lambda_p \|\bar{\mathbf{W}}^-\|_1 \right\}$$

subject to $\Omega = \bar{\mathbf{W}}$. The scaled augmented Lagrangian for this problem is [49]

$$\begin{aligned} \mathcal{L}_1^{aL}(\Omega, \bar{\mathbf{W}}, \bar{\mathbf{U}}) &= (1/p)(\text{tr}(\tilde{\Theta}\Omega) - \ln(|\Omega|)) \\ &+ \lambda_p \|\bar{\mathbf{W}}^-\|_1 + \frac{\rho}{2} \|\Omega - \bar{\mathbf{W}} + \bar{\mathbf{U}}\|_F^2 \end{aligned}$$

where $\bar{\mathbf{U}}$ is the dual variable, and $\rho > 0$ is the penalty parameter. Given the results $\tilde{\Omega}^{(i)}, \bar{\mathbf{W}}^{(i)}, \bar{\mathbf{U}}^{(i)}$ of the i th iteration, in the $(i+1)$ st iteration, an ADMM algorithm executes the following three updates until convergence:

Step (a). $\tilde{\Omega}^{(i+1)} \leftarrow \arg \min_{\Omega} \mathcal{L}_1^{aL}(\Omega, \bar{\mathbf{W}}^{(i)}, \bar{\mathbf{U}}^{(i)})$. We choose Ω to minimize

$$\text{tr}(\tilde{\Theta}\Omega) - \ln(|\Omega|) + \frac{p\rho}{2} \|\Omega - \bar{\mathbf{W}}^{(i)} + \bar{\mathbf{U}}^{(i)}\|_F^2.$$

The solution is as follows [40]. Let $\mathbf{Q}\mathbf{J}\mathbf{Q}^\top$ denote the eigen-decomposition of $\tilde{\Theta} - p\rho(\bar{\mathbf{W}}^{(i)} - \bar{\mathbf{U}}^{(i)})$ with diagonal matrix

\mathbf{J} consisting of the eigenvalues and $\mathbf{Q}\mathbf{Q}^\top = \mathbf{Q}^\top\mathbf{Q} = \mathbf{I}_q$. Then $\tilde{\Omega}^{(i+1)} = \mathbf{Q}\tilde{\mathbf{J}}\mathbf{Q}^\top$ where $\tilde{\mathbf{J}}$ is the diagonal matrix with ℓ th diagonal element

$$[\tilde{\mathbf{J}}]_{\ell\ell} = \frac{-[\mathbf{J}]_{\ell\ell} + \sqrt{([\mathbf{J}]_{\ell\ell})^2 + 4p\rho}}{2p\rho}.$$

Step (b). $\bar{\mathbf{W}}^{(i+1)} \leftarrow \arg \min_{\bar{\mathbf{W}}} \mathcal{L}_1^{aL}(\tilde{\Omega}^{(i+1)}, \bar{\mathbf{W}}, \bar{\mathbf{U}}^{(i)})$. We update $\bar{\mathbf{W}}^{(i+1)}$ as the minimizer w.r.t. $\bar{\mathbf{W}}$ of

$$\lambda_p \|\bar{\mathbf{W}}^-\|_1 + \frac{\rho}{2} \|\tilde{\Omega}^{(i+1)} - \bar{\mathbf{W}} + \bar{\mathbf{U}}^{(i)}\|_F^2.$$

The solution is soft thresholding given by [40]

$$[\bar{\mathbf{W}}]_{jk}^{(i+1)} = \begin{cases} [\tilde{\Omega}^{(i+1)} - \bar{\mathbf{U}}^{(i)}]_{jj} & \text{if } j = k \\ S_F([\tilde{\Omega}^{(i+1)} - \bar{\mathbf{U}}^{(i)}]_{jk}, \frac{\lambda_p}{\rho}) & \text{if } j \neq k \end{cases}$$

where $S_F(\cdot)$ denotes soft-thresholding as in Sec. IV-C.

Step (c). $\bar{\mathbf{U}}^{(i+1)} \leftarrow \bar{\mathbf{U}}^{(i)} + (\tilde{\Omega}^{(i+1)} - \bar{\mathbf{W}}^{(i+1)})$.

D. Practical Implementation

Here we present our implementation of the algorithms of Secs. IV-A-IV-C that was used in our numerical results.

- I. Parameters $\bar{\mu} = 10$, $\tau_{rel} = \tau_{abs} = 10^{-4}$, $\tau_{ff} = 10^{-5}$, $m_{\max} = 20$, $i_{\max} = 100$ and $\rho^{(0)} = 2$. Initialize $m = 1$, $\Omega^{(0)} = \mathbf{I}_p$, $\Phi_k^{(0)} = \mathbf{I}_q$, $k \in [M]$.
- II. For $m = 1, 2, \dots, m_{\max}$, do steps III-IV.
- III. Set $\hat{\Omega} = \Omega^{(m-1)}$ in (27). Pick $\rho = \rho^{(0)}$, $\tilde{\Phi}_k^{(0)} = \mathbf{I}_q$ for $k \in [M]$. For $i = 0, 1, \dots, i_{\max}$, do steps 1-6 below.

1. For $k \in [M]$, update Φ_k as $\tilde{\Phi}_k^{(i+1)}$ as in step (a), Sec. IV-B, then update \mathbf{W}_k as $\mathbf{W}_k^{(i+1)}$ as in step (b), Sec. IV-B, and then update \mathbf{U}_k as $\mathbf{U}_k^{(i+1)}$ as in step (c), Sec. IV-B, all with $\rho = \rho^{(i)}$.
2. Check for convergence following [9, Sec. 4.1.5]. Define the primal residual matrix $\mathbf{E}_{pri}^{(i+1)} \in \mathbb{C}^{q \times (qM)}$ at the $(i+1)$ st iteration as

$$\mathbf{E}_{pri}^{(i+1)} = \begin{bmatrix} \tilde{\Phi}_1^{(i+1)} - \mathbf{W}_1^{(i+1)}, \tilde{\Phi}_2^{(i+1)} - \mathbf{W}_2^{(i+1)}, \dots, \\ \tilde{\Phi}_M^{(i+1)} - \mathbf{W}_M^{(i+1)} \end{bmatrix}$$

and the dual residual matrix $\mathbf{E}_{dual}^{(i+1)} \in \mathbb{C}^{q \times (qM)}$ at the $(i+1)$ st iteration as

$$\mathbf{E}_{dual}^{(i+1)} = \rho^{(i)} \begin{bmatrix} \mathbf{W}_1^{(i+1)} - \mathbf{W}_1^{(i)}, \mathbf{W}_2^{(i+1)} - \mathbf{W}_2^{(i)}, \dots, \\ \mathbf{W}_M^{(i+1)} - \mathbf{W}_M^{(i)} \end{bmatrix}.$$

Let $e_1 = \|\tilde{\Phi}_1^{(i+1)} - \mathbf{W}_1^{(i+1)} \dots \tilde{\Phi}_M^{(i+1)} - \mathbf{W}_M^{(i+1)}\|_F$, $e_2 = \|\mathbf{W}_1^{(i+1)} - \mathbf{W}_1^{(i)} \dots \mathbf{W}_M^{(i+1)} - \mathbf{W}_M^{(i)}\|_F$, $e_3 = \|\mathbf{U}_1^{(i+1)} - \mathbf{U}_1^{(i)} \dots \mathbf{U}_M^{(i+1)} - \mathbf{U}_M^{(i)}\|_F$, $\tau_{pri} = q\sqrt{M}\tau_{abs} + \tau_{rel} \max(e_1, e_2)$ and $\tau_{dual} = q\sqrt{M}\tau_{abs} + \tau_{rel} e_3/\rho^{(i)}$. If $\|\mathbf{E}_{pri}^{(i+1)}\|_F \leq \tau_{pri}$ and $\|\mathbf{E}_{dual}^{(i+1)}\|_F \leq \tau_{dual}$, the convergence criterion is met. If the convergence criterion is met or if $i+1 > i_{\max}$, exit to step IV after setting $\Phi_k^{(m)} = \tilde{\Phi}_k^{(i+1)}$, $k \in [M]$, and then normalizing $\Gamma^{(m)} \leftarrow \Gamma^{(m)}/\|\Gamma^{(m)}\|_F$, else continue,

3. Update variable penalty parameter ρ as

$$\rho^{(i+1)} = \begin{cases} 2\rho^{(i)} & \text{if } \|\mathbf{E}_{pri}^{(i+1)}\|_F > \bar{\mu}\|\mathbf{E}_{dual}^{(i+1)}\|_F \\ \rho^{(i)}/2 & \text{if } \|\mathbf{E}_{dual}^{(i+1)}\|_F > \bar{\mu}\|\mathbf{E}_{pri}^{(i+1)}\|_F \\ \rho^{(i)} & \text{otherwise.} \end{cases}$$

For $k \in [M]$, set $\mathbf{U}_k^{(i+1)} = \mathbf{U}_k^{(i)}/2$ if $\|\mathbf{E}_{pri}^{(i+1)}\|_F > \bar{\mu}\|\mathbf{E}_{dual}^{(i+1)}\|_F$ and $\mathbf{U}_k^{(i+1)} = 2\mathbf{U}_k^{(i)}$ if $\|\mathbf{E}_{dual}^{(i+1)}\|_F > \bar{\mu}\|\mathbf{E}_{pri}^{(i+1)}\|_F$.

4. Set $i \leftarrow i + 1$ and return to step 2.

IV. Set $\hat{\Gamma} = \Gamma^{(m)}$ in (25). Pick $\rho = \rho^{(0)}$, $\tilde{\Omega}^{(0)} = \mathbf{I}_p$. For $i = 0, 1, \dots, i_{\max}$, do steps i-v below.

i. Update Ω as $\tilde{\Omega}^{(i+1)}$ as in step (a), Sec. IV-C, then update $\bar{\mathbf{W}}$ as $\bar{\mathbf{W}}^{(i+1)}$ as in step (b), Sec. IV-C, and then update $\bar{\mathbf{U}}$ as $\bar{\mathbf{U}}^{(i+1)}$ as in step (c), Sec. IV-C, all with $\rho = \rho^{(i)}$.

ii. Check for convergence following [40, Sec. II-A]. Define the primal residual matrix $\mathbf{H}_{pri}^{(i+1)} = \tilde{\Omega}^{(i+1)} - \bar{\mathbf{W}}^{(i+1)}$ and the dual residual matrix $\mathbf{H}_{dual}^{(i+1)} = \rho^{(i)}[\bar{\mathbf{W}}^{(i+1)} - \bar{\mathbf{W}}^{(i)}]$ where $\mathbf{H}_{pri}^{(i+1)}, \mathbf{H}_{dual}^{(i+1)} \in \mathbb{C}^{p \times p}$. Let

$$\begin{aligned} \tau_{pri} &= p\tau_{abs} + \tau_{rel} \max(\|\tilde{\Omega}^{(i+1)}\|_F, \|\bar{\mathbf{W}}^{(i+1)}\|_F) \\ \tau_{dual} &= p\tau_{abs} + \tau_{rel} \|\bar{\mathbf{U}}^{(i+1)}\|_F / \rho^{(i)}. \end{aligned}$$

If $\|\mathbf{H}_{pri}^{(i+1)}\|_F \leq \tau_{pri}$ and $\|\mathbf{H}_{dual}^{(i+1)}\|_F \leq \tau_{dual}$, the convergence criterion is met. If the convergence criterion is met or if $i + 1 > i_{\max}$, exit to step V after setting $\Omega^{(m)} = \tilde{\Omega}^{(i+1)}$, else continue.

iii. Update variable penalty parameter ρ as

$$\rho^{(i+1)} = \begin{cases} 2\rho^{(i)} & \text{if } \|\mathbf{H}_{pri}^{(i+1)}\|_F > \bar{\mu}\|\mathbf{H}_{dual}^{(i+1)}\|_F \\ \rho^{(i)}/2 & \text{if } \|\mathbf{H}_{dual}^{(i+1)}\|_F > \bar{\mu}\|\mathbf{H}_{pri}^{(i+1)}\|_F \\ \rho^{(i)} & \text{otherwise} \end{cases}$$

Set $\bar{\mathbf{U}}^{(i+1)} = \bar{\mathbf{U}}^{(i)}/2$ if $\|\mathbf{H}_{pri}^{(i+1)}\|_F > \bar{\mu}\|\mathbf{H}_{dual}^{(i+1)}\|_F$ and $\bar{\mathbf{U}}^{(i+1)} = 2\bar{\mathbf{U}}^{(i)}$ if $\|\mathbf{H}_{dual}^{(i+1)}\|_F > \bar{\mu}\|\mathbf{H}_{pri}^{(i+1)}\|_F$.

iv. Set $i \leftarrow i + 1$ and return to step ii.

V. Check for convergence of the flip-flop algorithm. If $\|\Gamma^{(m)} - \Gamma^{(m-1)}\|_F / \|\Gamma^{(m-1)}\|_F \leq \tau_{ff}$ and $\|\Omega^{(m)} - \Omega^{(m-1)}\|_F / \|\Omega^{(m-1)}\|_F \leq \tau_{ff}$, or $m > m_{\max}$, go to step VI, else set $m \leftarrow m + 1$ and return to step III.

VI. The final estimates are given by $\hat{\Omega} = \Omega^{(m)}$ and $\hat{\Gamma} = \Gamma^{(m)}$, and $\mathcal{E}_{\hat{\Omega}} = \{(i, j) : |[\hat{\Omega}]_{ij}| > 0\}$ and $\mathcal{E}_{\hat{\Gamma}} = \{(i, j) : \|\hat{\Phi}^{(ij)}\| > 0\}$ are the estimated edgesets for Ω and Γ respectively.

Remark 1. We terminate the flip-flop optimization (step V) when relative improvements in new updates of both $\Omega^{(m)}$ and $\Gamma^{(m)}$ are below the threshold τ_{ff} , or the maximum number of iterations in m is reached. The ADMM algorithms are terminated when both primary and dual residuals are below the respective tolerances, or the maximum number of iterations in i is reached; here we follow [49, Sec. 3.3.1] (see also [40] and [9]). The variable penalty $\rho^{(i)}$ follows the recommendations in [49, Sec. 3.4.1]. The most expensive computation in Sec. IV-B is in step (a) requiring the eigen-decomposition of M $q \times q$ matrices, with computational complexity $\mathcal{O}(Mq^3)$. Similarly, the most expensive computation in Sec. IV-C is in step (a) requiring the eigen-decomposition of a $p \times p$ matrix, with computational complexity $\mathcal{O}(p^3)$. Thus the overall computational complexity of our proposed approach is $\mathcal{O}(Mq^3 + p^3)$. \square

E. BIC for selection of λ_p , λ_q (and α)

Given n , K and M , the Bayesian information criterion (BIC) is given by (see also [9])

$$\begin{aligned} \text{BIC}(\lambda_p, \lambda_q, \alpha) &= -2KMq \ln(|\hat{\Omega}|) \\ &+ 2Kp \sum_{k=1}^M \left(-\ln(|\hat{\Phi}_k|) + p^{-1} \text{Re}(\text{tr}(\hat{\mathbf{A}}_k)) \right) \\ &+ \ln(2KM) \left(|\hat{\Omega}|_0 / 2 + \sum_{k=1}^M |\hat{\Phi}_k|_0 \right) \end{aligned} \quad (28)$$

where $\hat{\mathbf{A}}_k$ is given by (19) with Ω and Φ_k therein replaced with $\hat{\Omega}$ and $\hat{\Phi}_k$, respectively, $|\mathbf{H}|_0$ denotes number of nonzero elements in \mathbf{H} , $2KM$ is total number of real-valued measurements in frequency-domain and $2K$ is the number of real-valued measurements per frequency point, with total M frequencies in $(0, 0.5)$. A general expression for BIC is $-2 \log\text{-likelihood} + (\text{number of model parameters}) \times \log(\text{number of data points})$. The expression in (28) follows by using $\{\mathbf{d}_z(f_m)\}_{m=1}^{(n/2)-1}$ as complex-valued data in frequency-domain whose log-likelihood is given by (18). We count each complex value as two real values, both for data points and for parameters (entries of $\hat{\Phi}_k$), and also use the fact that $\hat{\Omega}$ is symmetric and $\hat{\Phi}_k$ is Hermitian.

Pick α , λ_q and λ_p to minimize BIC. In our simulations we fixed $\alpha = 0.05$ and then picked λ_q and λ_p over a grid of values, as follows. We search over $\lambda_q \in [\lambda_{q\ell}, \lambda_{qu}]$ and $\lambda_p \in [\lambda_{p\ell}, \lambda_{pu}]$ selected via a heuristic as in [40]. Find the smallest λ_q and λ_p , labeled λ_{qsm} and λ_{psm} , for which we get a no-edge model; then we set $\lambda_{qu} = \lambda_{qsm}/2$ and $\lambda_{q\ell} = \lambda_{qu}/10$; similarly for λ_{pu} and $\lambda_{p\ell}$.

V. THEORETICAL ANALYSIS

Now we provide sufficient conditions for local convergence in the Frobenius norm of the Kronecker-decomposable inverse PSD estimators to the true value or a scaled version of it. First some notation. The true values of Ω , Σ and $\bar{\mathbf{S}}(f)$ will be denoted as Ω° , Σ° and $\bar{\mathbf{S}}^\circ(f)$, respectively. Therefore, $\Omega^\circ = (\Sigma^\circ)^{-1}$. Since we use $\Phi_k := \bar{\mathbf{S}}^{-1}(f_k)$, we have $\Phi_k^\circ := \bar{\mathbf{S}}^{-\circ}(f_k)$ (where $\mathbf{A}^{-\circ} = (\mathbf{A}^\circ)^{-1}$). Therefore, in this notation, $\mathbf{d}_z(f_m) \sim \mathcal{N}_c(\mathbf{0}, \mathbf{S}_z^\circ(f_m))$ and $\mathbf{S}_z^\circ(f_m) = \bar{\mathbf{S}}^\circ(f_m) \otimes \Sigma^\circ$. Also in this section, we replace $\hat{\Phi}_k$'s in (25) with Φ_k 's and still use the notation $\hat{\Theta}$ for the sum (25) and the notation $\mathcal{L}_1(\Omega)$ for (24), and similarly, we replace $\hat{\Omega}$ in (27) with Ω and still use the notation $\hat{\Theta}$ for the sum (27) and $\mathcal{L}_2(\Gamma)$ for (26).

We follow [22] in first considering the solution to the unpenalized population objective function (i.e., expectation of $G(\Omega, \{\Phi\}, \{\Phi^*\})$ given by 18). We have

$$\begin{aligned} \bar{G}(\Omega, \{\Phi\}, \{\Phi^*\}) &= E\{G(\Omega, \{\Phi\}, \{\Phi^*\})\} \\ &= -\frac{1}{p} \ln(|\Omega|) - \frac{1}{2Mq} \sum_{k=1}^M \left[\ln(|\Phi_k|) + \ln(|\Phi_k^*|) \right] \\ &\quad - \frac{1}{p} \left(\text{tr}(\bar{\mathbf{S}}_k^\circ \Phi_k) + \text{tr}(\bar{\mathbf{S}}_k^\circ \Phi_k^*) \right) \text{tr}(\Sigma^\circ \Omega) \end{aligned} \quad (29)$$

where we have used the facts $\bar{\mathbf{S}}_k^\diamond = \bar{\mathbf{S}}^\diamond(\tilde{f}_k)$,

$$\begin{aligned} E\{\text{tr}(\mathbf{A}_k)\} &= \frac{1}{Kp} \sum_{\ell=-m_t}^{m_t} \text{tr}(E\{d_z(\tilde{f}_{k,\ell})d_z^H(\tilde{f}_{k,\ell})\}(\Phi_k \otimes \Omega)) \\ &= \frac{1}{p} \text{tr}((\bar{\mathbf{S}}_k^\diamond \otimes \Sigma^\diamond)(\Phi_k \otimes \Omega)) = \frac{1}{p} \text{tr}((\bar{\mathbf{S}}_k^\diamond \Phi_k) \otimes (\Sigma^\diamond \Omega)) \\ &= p^{-1} \text{tr}(\bar{\mathbf{S}}_k^\diamond \Phi_k) \text{tr}(\Sigma^\diamond \Omega). \end{aligned} \quad (30)$$

Define

$$\bar{\Omega}(\Gamma) = \arg \min_{\Omega} \bar{G}(\Omega, \{\Phi\}, \{\Phi^*\}), \quad (31)$$

$$\bar{\Gamma}(\Omega) = \arg \min_{\Gamma} \bar{G}(\Omega, \{\Phi\}, \{\Phi^*\}) \quad (32)$$

where $\bar{\Gamma}(\Omega) = [\bar{\Phi}_1(\Omega) \bar{\Phi}_2(\Omega) \cdots \bar{\Phi}_M(\Omega)]$.

Theorem 1. If $\sum_{k=1}^M (\text{tr}(\bar{\mathbf{S}}_k^\diamond \Phi_k) + \text{tr}(\bar{\mathbf{S}}_k^\diamond \Phi_k^*)) \neq 0$, then

$$\bar{\Omega}(\Gamma) = \frac{2Mq}{\sum_{k=1}^M (\text{tr}(\bar{\mathbf{S}}_k^\diamond \Phi_k) + \text{tr}(\bar{\mathbf{S}}_k^\diamond \Phi_k^*))} \Omega^\diamond, \quad (33)$$

and if $\text{tr}(\Sigma^\diamond \Omega) \neq 0$, then for $k \in [M]$,

$$\bar{\Phi}_k(\Omega) = \frac{p}{\text{tr}(\Sigma^\diamond \Omega)} (\bar{\mathbf{S}}_k^\diamond)^{-1} = \frac{p}{\text{tr}(\Sigma^\diamond \Omega)} \Phi_k^\diamond \bullet \quad (34)$$

Theorem 1 shows that the unpenalized population objective function yields true values up to a constant scalar. Notice that $\bar{\Omega}(\Gamma) = \Omega^\diamond$ if $\Gamma = \Gamma^\diamond$, and similarly, $\bar{\Phi}_k(\Omega) = \Phi_k^\diamond$, $k = 1, 2, \dots, M$, if $\Omega = \Omega^\diamond$.

We now turn to penalized data-based objective function $\mathcal{L}(\Omega, \Gamma)$ which is minimized alternately as $\mathcal{L}_2(\Gamma)$ w.r.t. Φ_k 's and as $\mathcal{L}_1(\Omega)$ w.r.t. Ω . Here in addition to assumptions (A1) and (A2), we assume

(A3) Define the true edgesets $\mathcal{S}_q = \{\{i, j\} : [(\bar{\mathbf{S}}^\diamond)^{-1}(f)]_{ij} \neq 0, i \neq j, 0 \leq f \leq 0.5, i, j \in [q]\}$ and $\mathcal{S}_p = \{\{i, j\} : [\Omega^\diamond]_{ij} \neq 0, i \neq j, i, j \in [p]\}$, where $\bar{\mathbf{S}}^\diamond(f)$ denotes DTFT of $\Psi(\tau)$ and $\Omega^\diamond = (\Sigma^\diamond)^{-1}$ denotes the true value of Ω . Assume that number of nonzero elements in the true edgesets \mathcal{S}_q and \mathcal{S}_p are upperbounded as $|\mathcal{S}_q| \leq s_q$ and $|\mathcal{S}_p| \leq s_p$.

(A4) The minimum and maximum eigenvalues of $q \times q$ PSD $\bar{\mathbf{S}}^\diamond(f) \succ \mathbf{0}$ satisfy $0 < \beta_{q,\min} \leq \min_{f \in [0,0.5]} \phi_{\min}(\bar{\mathbf{S}}^\diamond(f))$ and $\max_{f \in [0,0.5]} \phi_{\max}(\bar{\mathbf{S}}^\diamond(f)) \leq \beta_{q,\max} < \infty$. Similarly, $0 < \beta_{p,\min} \leq \phi_{\min}(\Sigma^\diamond) \leq \phi_{\max}(\Sigma^\diamond) \leq \beta_{p,\max} < \infty$. Here $\beta_{\cdot,\min}$ and $\beta_{\cdot,\max}$ are not functions of n, p, q .

Theorem 2 establishes bounds on estimation errors of local minimizers $\hat{\Omega}(\Gamma)$ and $\hat{\Gamma}(\Omega)$ of $\mathcal{L}_1(\Omega)$ and $\mathcal{L}_2(\Gamma)$, respectively. We now explicitly allow $p, q, M, K, s_p, s_q, \lambda_p$ and λ_q to be functions of sample size n , denoted as $p_n, q_n, M_n, K_n, s_{p_n}, s_{q_n}, \lambda_{p_n}$ and λ_{q_n} , respectively. (In the appendices we do not do so to keep the notation simple.) First we define some variables. For $\tau > 2$, define

$$\gamma_p = 0.1/\beta_{p,\max}, \quad (35)$$

$$C_{1p} = \frac{2}{\sqrt{\ln(M_n^{1/\tau} q_n)}} + \sqrt{2\tau + \frac{2 \ln(16)}{\ln(M_n^{1/\tau} q_n)}}, \quad (36)$$

$$C_{0q} = 16C_{1q}(1 + \gamma_p \beta_{p,\max}) \beta_{q,\max}, \quad (37)$$

$$\gamma_q = 0.1/\beta_{q,\max}, \quad (38)$$

$$C_{1p} = \sqrt{\frac{2}{\ln(p_n)}} + \sqrt{\tau + \frac{\ln(4)}{\ln(p_n)}}, \quad (39)$$

$$C_{0p} = 8C_{1p}(2 + \gamma_q \beta_{q,\max}) \beta_{p,\max}, \quad (40)$$

$$r_{qn} = \sqrt{M_n(q_n + s_{q_n}) \ln(M_n^{1/\tau} q_n)/(K_n p_n)} = o(1), \quad (41)$$

$$r_{pn} = \sqrt{(p_n + s_{p_n}) \ln(p_n)/(M_n K_n q_n)} = o(1). \quad (42)$$

Theorem 2. Let $\tau > 2$.

(i) Let $\mathcal{B}(\Gamma^\diamond) = \{\Gamma : \|\Gamma - \Gamma^\diamond\|_F \leq \gamma_q, \Phi_k = \Phi_k^H \succ \mathbf{0}\}$ and $\hat{\Omega}(\Gamma) = \arg \min_{\{\Omega: \Gamma \in \mathcal{B}(\Gamma^\diamond)\}} \mathcal{L}_1(\Omega)$. Suppose λ_{p_n} satisfies

$$\frac{C_{0p}}{p_n} \sqrt{\frac{\ln(p_n)}{M_n K_n q_n}} \leq \lambda_{p_n} \leq \frac{C_{0p}}{p_n} \sqrt{\left(1 + \frac{p_n}{s_{p_n}}\right) \frac{\ln(p_n)}{M_n K_n q_n}}. \quad (43)$$

Let $N_p := \arg \min_n \{n : r_{pn} \leq \beta_{p,\min}/(34C_{0p})\}$. Then under assumptions (A1)-(A4), for $n > N_p$, $\hat{\Omega}(\Gamma)$ satisfies

$$\|\hat{\Omega}(\Gamma) - \bar{\Omega}(\Gamma)\|_F \leq \frac{17C_{0p}}{\beta_{p,\min}^2} r_{pn} \quad (44)$$

with probability greater than $1 - \frac{1}{p^{\tau-2}} - 4p^2 e^{-KqM}$.

(ii) Let $\mathcal{B}(\Omega^\diamond) = \{\Omega : \|\Omega - \Omega^\diamond\|_F \leq \gamma_p, \Omega = \Omega^\top \succ \mathbf{0}\}$ and $\hat{\Gamma}(\Omega) = \arg \min_{\{\Gamma: \Omega \in \mathcal{B}(\Omega^\diamond)\}} \mathcal{L}_2(\Gamma)$. Suppose λ_{q_n} satisfies

$$\frac{C_{0q}}{M_n q_n} \sqrt{d_n} \leq \lambda_{q_n} \leq \frac{C_{0q}}{M_n q_n} \sqrt{\left(1 + \frac{q_n}{s_{q_n}}\right) d_n}, \quad (45)$$

where $d_n = \ln(M_n^{1/\tau} q_n)/(K_n p_n)$. Let $N_q := \arg \min_n \{n : r_{qn} \leq \beta_{q,\min}/(34C_{0q})\}$. Then under assumptions (A1)-(A4), for $n > N_q$ and $\alpha \in [0, 1]$, $\hat{\Gamma}(\Omega)$ satisfies

$$\|\hat{\Gamma}(\Omega) - \bar{\Gamma}(\Omega)\|_F \leq \frac{17C_{0q}}{\beta_{q,\min}^2} r_{qn} \quad (46)$$

with probability greater than $1 - \frac{1}{q^{\tau-2}} - 16Mq^2 e^{-Kp/2}$ •

Remark 2. Theorem 2 helps determine how to choose M_n and K_n so that for given n, s_{p_n}, s_{q_n}, q_n and p_n , $\lim_{n \rightarrow \infty} r_{pn} = 0$ and $\lim_{n \rightarrow \infty} r_{qn} = 0$, and moreover, how fast can s_{p_n} and s_{q_n} grow with n and still have r_{qn} and $r_{pn} \downarrow 0$. Since $K_n M_n \approx n/2$, if one picks $K_n = \mathcal{O}(n^\mu)$, then $M_n = \mathcal{O}(n^{1-\mu})$ for some $0 < \mu < 1$. We assume p_n and q_n are of the same order. (i) First consider the case $\mathcal{O}(p_n) = \mathcal{O}(p_n + s_{p_n}) = \mathcal{O}(q_n) = \mathcal{O}(q_n + s_{q_n})$, which, for example, is true for chain graphs. Also, take $\mathcal{O}(p_n) \propto n^\nu$ for some $\nu > 0$. Then $r_{pn} = \mathcal{O}(\sqrt{\ln(n)/n}) \rightarrow 0$ as $n \rightarrow \infty$, and $r_{qn} = \mathcal{O}(\sqrt{\ln(n)/n^{2\mu-1}}) \rightarrow 0$ as $n \rightarrow \infty$ if $\mu > 0.5$. This holds for any $\nu > 0$. If $\mu = \frac{3}{4}$, then $r_{qn} = \mathcal{O}(\sqrt{\ln(n)/n^{1/4}}) > r_{pn}$. If $\mu = \frac{2}{3}$, then $r_{qn} = \mathcal{O}(\sqrt{\ln(n)/n^{1/6}}) > r_{pn}$. (ii) Now suppose $\mathcal{O}(p_n) = \mathcal{O}(q_n) \propto n^\nu$ for some $\nu > 0$, but $\mathcal{O}(s_{p_n}) = \mathcal{O}(s_{q_n}) \propto n^{2\nu} = \mathcal{O}(p_n^2)$, which is true for Erdős-Rényi graphs, e.g. Then $r_{pn} = \mathcal{O}(\sqrt{\ln(n)/n^{1-\nu}}) \rightarrow 0$ as $n \rightarrow \infty$ if $\nu < 1$, and $r_{qn} = \mathcal{O}(\sqrt{\ln(n)/n^{2\mu-1-\nu}}) \rightarrow 0$ as $n \rightarrow \infty$ if $2\mu - \nu > 1$. Clearly $\nu = 1$ does not work. Suppose $\nu = 0.25$ and $\mu = 0.75$. Then $r_{pn} = \mathcal{O}(\sqrt{\ln(n)/n^{0.375}})$ and $r_{qn} = \mathcal{O}(\sqrt{\ln(n)/n^{1/8}})$. □

Remark 3. The values of γ_p and γ_q specified in (35) and (38), respectively, are used in the proofs of Theorem 2(ii) (see after (116)) and Theorem 2(i) (see (101)), respectively. One can enlarge γ_p and γ_q to $\gamma_p = 0.1\sqrt{p_n}/\beta_{p,\min}$ and $\gamma_q = 0.1\sqrt{M_n q_n}/\beta_{q,\min}$, respectively, and the proofs and the other results remain unchanged and valid. Enlarging these values implies that the balls $\mathcal{B}(\Gamma^\circ)$ and $\mathcal{B}(\Omega^\circ)$ specified in Theorem 2 are larger, signifying larger convergence regions for initialization of Γ and Ω . However, this would slow the convergence rates from $\|\hat{\Omega}(\Gamma) - \bar{\Omega}(\Gamma)\|_F = \mathcal{O}_P(r_{pn})$ and $\|\hat{\Gamma}(\Omega) - \bar{\Gamma}(\Omega)\|_F = \mathcal{O}_P(r_{qn})$ to $\|\hat{\Omega}(\Gamma) - \bar{\Omega}(\Gamma)\|_F = \mathcal{O}_P(\sqrt{p_n} r_{pn})$ and $\|\hat{\Gamma}(\Omega) - \bar{\Gamma}(\Omega)\|_F = \mathcal{O}_P(\sqrt{M_n q_n} r_{qn})$, respectively. \square

Theorem 3. Assume $\|\Omega^\circ\|_F = 1$.

(i) Define $\hat{\Omega} = \hat{\Omega}(\Gamma)/\|\hat{\Omega}(\Gamma)\|_F$. Let $N_{2p} := \arg \min_n \{n : r_{pn} \leq \beta_{p,\min}^2 \|\hat{\Omega}(\Gamma)\|_F / (34C_{0p})\}$ and $\gamma_r = (\beta_{q,\max} + \beta_{q,\min})/\beta_{q,\min}$. Under the assumptions of Theorem 2(i), for $n > \max\{N_p, N_{2p}\}$, $\hat{\Omega}$ satisfies

$$\|\hat{\Omega} - \Omega^\circ\|_F \leq 4\gamma_r \|\hat{\Omega}(\Gamma) - \bar{\Omega}(\Gamma)\|_F \leq \frac{68\gamma_r C_{0p}}{\beta_{p,\min}^2} r_{pn} \quad (47)$$

with probability greater than $1 - \frac{1}{p^{\tau-2}} - 4p^2 e^{-KqM}$.

(ii) Let $C_{2p} = 68\gamma_r \sqrt{p}\beta_{p,\max} C_{0p} / \beta_{p,\min}^2$, $U_{1p} = p/(2C_{2p})$, $U_{2p} = 0.1\beta_{p,\min}^2 / (68\gamma_r C_{0p} \beta_{p,\max})$, $N_{3p} := \arg \min_n \{n : r_{pn} \leq \max\{U_{1p}, U_{2p}\}\}$ and $C_{2q} = 2C_{2p} \|\Gamma^\circ\|_F / p$. Let $\hat{\Gamma}(\hat{\Omega}) = \arg \min_{\{\Gamma: \Omega = \hat{\Omega} \in \mathcal{B}(\Omega^\circ)\}} \mathcal{L}_2(\Gamma)$ where $\hat{\Omega}$ is as in Theorem 3(i). Under the assumptions of Theorem 2, for $n > \max\{N_p, N_{2p}, N_q, N_{3p}\}$ and $\alpha \in [0, 1]$, $\hat{\Gamma}(\hat{\Omega})$ satisfies

$$\|\hat{\Gamma}(\hat{\Omega}) - \Gamma^\circ\|_F \leq \frac{17C_{0q}}{\beta_{q,\min}^2} r_{qn} + C_{2q} r_{pn} \quad (48)$$

with probability greater than $1 - \frac{1}{p^{\tau-2}} - 4p^2 e^{-KqM} - \frac{1}{q^{\tau-2}} - 16Mq^2 e^{-Kp/2}$ \bullet

VI. NUMERICAL RESULTS

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 linear conditional dependency structures underlying the data.

A. Synthetic Data

We use model (5)-(6) to generate synthetic data where $\Psi(\tau)$ is controlled via a vector autoregressive (VAR) model impulse response and Σ is determined via an Erdős-Rényi graph. We take $p = q = 15$. Consider the impulse response $\mathbf{H}_i^{(r)} \in \mathbb{R}^{5 \times 5}$ generated as $\mathbf{H}_i^{(r)} = \sum_{k=1}^3 \mathbf{A}_k^{(r)} \mathbf{H}_{i-k}^{(r)} + \mathbf{I}_5 \delta_i$, where $\mathbf{H}_i^{(r)} = 0$ for $i < 0$, δ_i is the Kronecker delta, $r = 1, 2, 3$, and only 5% of entries of $\mathbf{A}_i^{(r)}$'s are nonzero and the nonzero elements are independently and uniformly distributed over $[-0.8, 0.8]$. We then check if the VAR(3) model is stable with all eigenvalues of the companion matrix ≤ 0.95 in magnitude; if not, we re-draw randomly till this condition is fulfilled. The impulse response $\mathbf{B}_i \in \mathbb{R}^{15 \times 15}$ in (5) is given by $\mathbf{B}_i = \text{block-diag}\{\mathbf{H}_i^{(1)}, \mathbf{H}_i^{(2)}, \mathbf{H}_i^{(3)}\}$, for

$0 \leq i \leq L = 40$, otherwise it is set to zero. Thus \mathbf{B}_i 's in (5) have a block-diagonal structure with 3 blocks, each block is 5×5 . In the Erdős-Rényi graph with $p = 15$ nodes, the nodes are connected with probability $p_{er} = 0.05$. In the upper triangular $\bar{\Omega}$, $\bar{\Omega}_{ij} = 0$ if $\{i, j\} \notin \mathcal{S}_p$, $\bar{\Omega}_{ij}$ is uniformly distributed over $[-0.4, -0.1] \cup [0.1, 0.4]$ if $\{i, j\} \in \mathcal{S}_p$, and $\bar{\Omega}_{ii} = 0.5$. With $\bar{\Omega} = \bar{\Omega}^\top$, add $\kappa \mathbf{I}_p$ to $\bar{\Omega}$ with κ picked to make minimum eigenvalue of $\Omega = \bar{\Omega} + \kappa \mathbf{I}_p$ equal to 0.5. Let $\Omega = \tilde{\mathbf{F}} \tilde{\mathbf{F}}$ (matrix square-root), then $\mathbf{F} = \tilde{\mathbf{F}}^{-1}$ in (5).

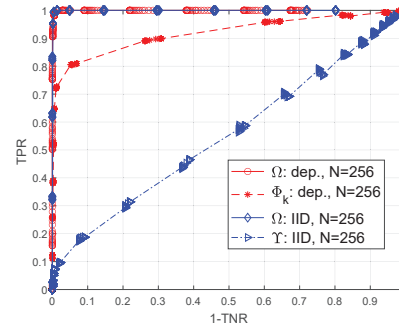


Fig. 1: ROC curves: plots labeled “IID” are from the approach of [17], [18], [25], and the plots labeled “dep.” are from our proposed approach. TPR=true positive rate, TNR=true negative rate

We applied our proposed approach with $n = 256$, $M = 2$, $K = 63$ and compared with the approach of [17] (which is also the approach of [18], [25], all of whom assume i.i.d. observations and have two lasso penalties one each on Ω and Υ , counterpart to our Φ_k). In our approach, we fix $\alpha = 0.05$ for all simulations and real data results. For fixed values of λ_q and λ_p , using our proposed approach of Sec. IV-D, we calculated the true positive rate (TPR) and false positive rate 1-TNR (where TNR is the true negative rate) over 100 runs, separately for Ω and $\{\Phi_k\}/\Upsilon$, based on the estimated edges. As we vary λ_q and λ_p over a wide range of values, we can compute the corresponding pairs of estimated (1-TNR, TPR). The receiver operating characteristic (ROC) is shown in Fig. 1 based on 100 runs, using the estimated (1-TNR, TPR). We repeat this method for the i.i.d. modeling approach of [17], [18], [25]. Fig. 1 shows that the i.i.d. modeling of [17], [18], [25] is unable to capture the “dependent” edges (cf. (4)) via Υ whereas it has no issues with Ω . Our approach works well for both components of the Kronecker product graph.

In Table I we show the results based on 100 runs under different parameter settings and samples sizes. Here we show the F_1 score, TPR, 1-TNR and timing values for the overall graph (not the two Kronecker product components separately) along with the $\pm\sigma$ errors. 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. We take $n = 64, 128, 256, 512, 1024, 2048$, and for our proposed approach, the corresponding m_t values leading to different M values are $m_t = 7, 15, 31, 63, 127, 255$ ($M = 2$), $m_t = 4, 9, 20, 41, 84, 169$ ($M = 3$), $m_t = 3, 7, 14, 31, 63, 127$ ($M = 4$), $m_t = 2, 5, 12, 24, 50, 101$ ($M = 5$), $m_t = **, 4, 10, 20, 42, 84$ ($M = 6$), $m_t = **, **, **, 15, 31, 63$ ($M = 8$), and $m_t = **, **, **, 12, 25, 50$ ($M = 10$). Here ** denotes that no simulation was performed for the corresponding sample size n (since $K = 2m_t + 1$ is too

TABLE I: F_1 scores, TPR, 1-TNR and timing per run for fixed tuning parameters, for the synthetic data example, averaged over 100 runs. The entries ** denote no simulations done for these parameters.

n	64	128	256	512	1024	2048
Proposed Approach: F_1 scores $\pm\sigma$ when λ 's are selected to minimize BIC						
$M=2$	0.5163 \pm 0.1530	0.5660 \pm 0.1580	0.6440 \pm 0.1709	0.7061 \pm 0.1147	0.7018 \pm 0.1176	0.7190 \pm 0.1217
$M=3$	0.5111 \pm 0.1705	0.5876 \pm 0.1560	0.6969 \pm 0.1421	0.7266 \pm 0.1159	0.7322 \pm 0.1031	0.7474 \pm 0.1000
$M=4$	0.5454 \pm 0.1852	0.6470 \pm 0.1489	0.7106 \pm 0.1465	0.7376 \pm 0.1011	0.7446 \pm 0.1069	0.7524 \pm 0.1047
$M=5$	0.5977 \pm 0.1717	0.6609 \pm 0.1465	0.7049 \pm 0.1367	0.7253 \pm 0.1104	0.7401 \pm 0.1028	0.7467 \pm 0.1002
$M=6$	**	0.6277 \pm 0.1353	0.6773 \pm 0.1379	0.7115 \pm 0.1172	0.7343 \pm 0.1025	0.7369 \pm 0.0971
$M=8$	**	**	**	0.7016 \pm 0.1071	0.7366 \pm 0.1013	0.7365 \pm 0.0974
$M=10$	**	**	**	0.7123 \pm 0.1117	0.7319 \pm 0.1044	0.7367 \pm 0.1022
Proposed Approach: F_1 scores $\pm\sigma$ when λ 's are selected to maximize F_1 score						
$M=2$	0.6826 \pm 0.1440	0.6954 \pm 0.1588	0.7485 \pm 0.1632	0.8026 \pm 0.1139	0.8032 \pm 0.1588	0.8440 \pm 0.1184
$M=3$	0.6984 \pm 0.1383	0.7322 \pm 0.1730	0.8055 \pm 0.1383	0.8293 \pm 0.1190	0.8372 \pm 0.1442	0.8670 \pm 0.1295
$M=4$	0.7041 \pm 0.1355	0.7364 \pm 0.1646	0.8074 \pm 0.1434	0.8282 \pm 0.1169	0.8401 \pm 0.1197	0.8633 \pm 0.1397
$M=5$	0.6652 \pm 0.1664	0.7309 \pm 0.1431	0.8072 \pm 0.1466	0.8411 \pm 0.1158	0.8451 \pm 0.1251	0.8637 \pm 0.1314
$M=6$	**	0.7218 \pm 0.1490	0.8089 \pm 0.1324	0.8252 \pm 0.1206	0.8433 \pm 0.1282	0.8583 \pm 0.1396
$M=8$	**	**	**	0.8329 \pm 0.1130	0.8382 \pm 0.1221	0.8601 \pm 0.1404
$M=10$	**	**	**	0.8187 \pm 0.1216	0.8286 \pm 0.1525	0.8496 \pm 0.1406
IID modeling [17], [18], [25]: λ 's are selected to maximize F_1 score						
F_1 scores $\pm\sigma$	0.4329 \pm 0.1244	0.4230 \pm 0.1208	0.4368 \pm 0.1228	0.4746 \pm 0.1367	0.4483 \pm 0.1180	0.4709 \pm 0.1104
timing (s) per run $\pm\sigma$	0.0051 \pm 0.0011	0.0073 \pm 0.0014	0.0111 \pm 0.0020	0.0195 \pm 0.0031	0.0342 \pm 0.0035	0.0640 \pm 0.0050
Proposed approach under model mismatch – non-Gaussian e in (5): F_1 scores $\pm\sigma$ when λ 's are selected to minimize BIC						
Exponential e , $M=4$	0.5518 \pm 0.1853	0.6565 \pm 0.1728	0.7098 \pm 0.1349	0.7355 \pm 0.0976	0.7514 \pm 0.1141	0.7555 \pm 0.0888
Uniform e , $M=4$	0.5434 \pm 0.1772	0.6510 \pm 0.1693	0.7137 \pm 0.1364	0.7400 \pm 0.1043	0.7494 \pm 0.1146	0.7537 \pm 0.0982
Proposed Approach: TPR $\pm\sigma$ when λ 's are selected to maximize F_1 score						
$M=2$	0.6312 \pm 0.1675	0.6420 \pm 0.1541	0.6937 \pm 0.1852	0.7533 \pm 0.1332	0.8146 \pm 0.1187	0.8249 \pm 0.1199
$M=4$	0.6793 \pm 0.1493	0.7120 \pm 0.1477	0.7595 \pm 0.1529	0.7919 \pm 0.1307	0.8142 \pm 0.1229	0.8836 \pm 0.1021
$M=6$	**	0.6711 \pm 0.1529	0.7459 \pm 0.1608	0.8024 \pm 0.1287	0.8162 \pm 0.1215	0.8275 \pm 0.1290
$M=10$	**	**	**	0.7867 \pm 0.1269	0.8278 \pm 0.1199	0.8504 \pm 0.1177
Proposed Approach: 1-TNR $\pm\sigma$ when λ 's are selected to maximize F_1 score						
$M=2$	0.0032 \pm 0.0092	0.0033 \pm 0.0090	0.0022 \pm 0.0074	0.0018 \pm 0.0061	0.0049 \pm 0.0157	0.0025 \pm 0.0096
$M=4$	0.0041 \pm 0.0118	0.0044 \pm 0.0127	0.0020 \pm 0.0074	0.0021 \pm 0.0097	0.0023 \pm 0.0086	0.0043 \pm 0.0174
$M=6$	**	0.0035 \pm 0.0116	0.0013 \pm 0.0050	0.0026 \pm 0.0113	0.0027 \pm 0.0120	0.0030 \pm 0.0161
$M=10$	**	**	**	0.0025 \pm 0.0113	0.0046 \pm 0.0173	0.0040 \pm 0.0173
Proposed Approach: timing (s) per run $\pm\sigma$ when λ 's are selected to minimize BIC						
$M=2$	0.1687 \pm 0.0400	0.1688 \pm 0.0418	0.1791 \pm 0.1005	0.1777 \pm 0.0289	0.2166 \pm 0.0322	0.3131 \pm 0.1051
$M=4$	0.2294 \pm 0.1650	0.2470 \pm 0.1026	0.2284 \pm 0.0846	0.2278 \pm 0.0392	0.2890 \pm 0.1338	0.3627 \pm 0.0494
$M=6$	**	0.2738 \pm 0.0903	0.2426 \pm 0.0633	0.2507 \pm 0.0537	0.2944 \pm 0.0902	0.3842 \pm 0.0471
$M=10$	**	**	**	0.3040 \pm 0.1400	0.3230 \pm 0.0733	0.4285 \pm 0.0890

small). We show the resulting F_1 scores under two different scenarios: when we use the proposed BIC parameter selection method (Sec. IV-E) and when F_1 score was selected based on λ values that maximize the F_1 score. While the latter approach is not practical, it is presented to illustrate what is possible using the proposed approach and what may be "lost" when there are errors in the BIC parameter selection method. The number of unknown parameters being estimated are $\mathcal{O}(p^2 + Mq^2)$, with $\mathcal{O}(p^2)$ for Ω and $\mathcal{O}(Mq^2)$ for $M \Phi_k$'s. We see that for a fixed n , at first the performance improves with increasing M , then it slowly declines as more parameters need to be estimated with increasing M . Increasing M also reduces $K = 2m_t + 1$ since $KM \approx \frac{n}{2}$, which reduces the number of frequency-domain samples (K) for averaging for the k th model for Φ_k , $k \in [M]$ (see assumption (A1) in Sec. III). Note also that by (46) of Theorem 2(ii), the error in estimating Φ_k 's $\propto r_{qn} \propto \sqrt{(Mq)/(Kp)}$. For a fixed M , the performance improves, in general, with increasing n but more slowly for higher n 's. Higher n values implies higher resolution in the frequency-domain and for fixed M , higher n implies higher K (and m_t), in which case assumption (A1) in Sec. III may not hold. The TPR, 1-TNR and timing values are shown for selected M 's for the proposed approach where timing per run is for the λ values picked by the BIC criterion. It is seen that increasing M and/or n leads to only a small

increase in timing.

In Table I we also show the performance of i.i.d. modeling approach of [17], [18], [25], in terms of the F_1 score and timing. The i.i.d. modeling approach is significantly faster but the accuracy in edge detection in terms of the F_1 score is much poorer. Finally, to assess sensitivity to modeling errors such as violation of the Gaussianity assumption, we used either exponential or uniform $e(t)$ in (5), both with zero-mean unit variance, instead of the assumed Gaussian $e(t)$ in our model. The results are shown for $M = 4$ and we see that the performance is robust w.r.t. violation of this assumption.

B. Real Data: Beijing air-quality dataset [50]

Here we consider Beijing air-quality dataset [50], [51], downloaded from <https://archive.ics.uci.edu/dataset/501/beijing+multi+site+air+quality+data>. This data set includes hourly air pollutants data from 12 nationally-controlled air-quality monitoring sites in the Beijing area. The time period is from March 1st, 2013 to February 28th, 2017. The six air pollutants are $PM_{2.5}$, PM_{10} , SO_2 , NO_2 , CO , and O_3 , 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 ($= q$) features (pollutants and weather variables). We used data from 8 ($= p$) sites: 4 rural/suburban sites Changping, Dingling, Huairou,

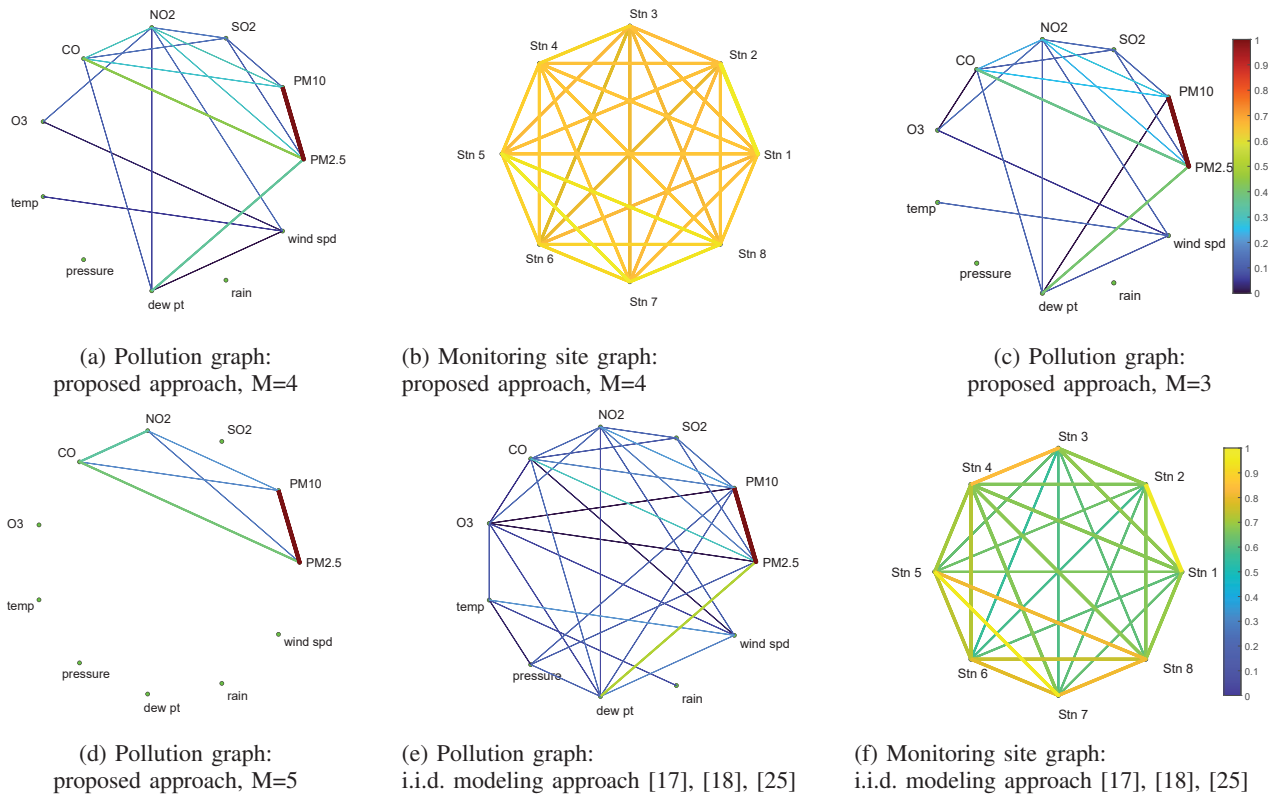


Fig. 2: Pollution and site graphs for the Beijing air-quality dataset [50] for year 2013-14: 8 monitoring sites and 11 features ($p = 8, q = 11, n = 364$). Number of distinct edges = 18, 28, 20, 6, 30, 28 in graphs (a), (b), (c), (d), (e) and (f), respectively. Monitoring sites labeled Stn. 1-4 are the rural/suburban sites and those labeled Stn. 5-8 are the urban sites (see the text). For the pollution graph, estimated $\hat{\Phi}^{(ij)}$ is the edge weight (normalized to have $\max_{i \neq j} \hat{\Phi}^{(ij)} = 1$) and for the site graph, estimated $|\hat{\Omega}_{ij}|$ is the edge weight (normalized to have $\max_{i \neq j} |\hat{\Omega}_{ij}| = 1$). The edge weights are color coded (all pollution graphs share the same color legend, and similarly for the site graphs), in addition to the edges with higher weights being drawn thicker.

Shunyi, and 4 urban sites Aotizhongxin, Dongsì, Guanyuan, Gucheng (labeled Stn 1 through 8 in Fig. 2). The data are averaged over 24 hour period to yield daily averages. We used one year 2013-14 of daily data resulting in $n = 365$ days. Arranging stations as rows and features as columns, we have $\mathbf{Z}(t) \in \mathbb{R}^{8 \times 11}, t = 1, 2, \dots, 365$. We pre-processed the data as follows. Given j th feature data $\mathbf{Z}_{ij}(t)$ at i th station, we transform it to $\bar{\mathbf{Z}}_{ij}(t) = \ln(\mathbf{Z}_{ij}(t)/\mathbf{Z}_{ij}(t-1))$ for each i and j , and then detrend it (i.e., remove the best straight-line fit). Finally, we scale the detrended scalar sequence to have a mean-square value of one. All temperatures were converted from Celsius to Kelvin to avoid negative numbers. If a value of a feature is zero (e.g., wind speed), we added a small positive number to it so that the log transformation is well-defined.

We applied our proposed approach with $M = 4, K = 45$ and $n = 364$ ($p = 8, q = 11$) and compared it with the i.i.d. modeling approach of [17], [18], [25]. The objective here is to compare the two approaches in estimation of the pollution (feature) graph and the site graph. The spatio-temporal data has a matrix structure and one is interested in learning two aspects of conditional dependencies: the relationship among the features via the pollution graph and the relationship among the sites via the site graph. We have not yet tested if our model assumptions apply to this dataset (this needs further theoretical analysis to devise suitable statistical tests, particularly in a high-dimensional setting), but it still seems to be useful to compare the results of our proposed approach and that of

[17], [18], [25]. Fig. 2(a) shows the resulting graph for the air quality and environmental variables where $\{i, j\} \in \mathcal{S}_q$ iff $\hat{\Phi}^{(ij)} = (\sum_{k=1}^M |[\hat{\Phi}_k]_{ij}|^2)^{1/2} > 0$ for $i \neq j$, and Fig. 2(b) shows the resulting graph for the sites around the Beijing area where $\{i, j\} \in \mathcal{S}_p$ iff $|\hat{\Omega}_{ij}| > 0$ for $i \neq j$. Since all the sites are physically close to one another, it is not surprising that the site graph in Fig. 2(b) is fully connected. But we do see that the rural/suburban sites stn. 1 through stn. 4 have higher weight edges among the group and the urban sites stn. 5 through stn. 8 have higher weight edges among the urban group, with inter-group edge weights being slightly weaker (but fully connected). Automobile exhaust is the main cause of NO_2 which is likely to undergo a chemical reaction with Ozone O_3 , thereby, lowering its concentration [51]. This fact is captured by the edge between NO_2 and Ozone O_3 in Fig. 2(a). Cold, dry air from the north of Beijing reduces both dew point and $\text{PM}_{2.5}$ particle concentration in suburban areas while southerly wind brings warmer and more humid air from the more polluted south that elevates both dew point and $\text{PM}_{2.5}$ concentration [50]. This fact is captured by the edge between dew point and $\text{PM}_{2.5}$ in Fig. 2(a). The counterparts to Figs. 2(a) and 2(b) when using the i.i.d. modeling approach of [17], [18], [25], are shown in Figs. 2(e) and 2(f), respectively. While the site graph in Fig. 2(f) is fully connected and quite similar to the proposed approach's site graph in Fig. 2(b), the pollution graph in Fig. 2(e) far denser than the proposed approach's

pollution graph in Fig. 2(a).

We do not have any systematic approach for selection of M for a given sample size n . Since $KM \approx \frac{n}{2}$, fixing M fixes $K = 2m_t + 1$, and vice-versa. Using BIC to pick M does not work as BIC always picks the smallest M . The synthetic data results presented in Table I show that the performance is not unduly sensitive to the choice of M . To illustrate the sensitivity of the proposed approach in Beijing data case, we show the pollution graphs in Figs. 2(c) and 2(d) for the choice $M = 3$ ($K=59$) and $M = 5$ ($K = 35$), respectively. There is not much difference between pollution graphs for $M = 4$ and $M = 3$, but that for $M = 5$ is much sparser. This is consistent with the results of Sec. VI-A on synthetic data.

VII. CONCLUSIONS

Sparse-group lasso penalized log-likelihood approach in frequency-domain with a Kronecker-decomposable PSD was investigated for matrix CIG learning for dependent time series. An ADMM-based flip-flop approach for iterative optimization of the bi-convex problem was presented. We provided sufficient conditions for consistency of a local estimator of inverse PSD. We illustrated our approach using numerical examples utilizing both synthetic and real data. Lasso and related approaches are known to yield biased estimates [52]. To remedy this, various non-convex penalties have been suggested [52] and typically, lasso-based approaches provide the initial guess for iterative optimization. In the context of this paper, adaptive lasso has been used in [40] (the basis of the ADMM method of Sec. IV-C), and a log-sum penalty has been used in [53] (which modifies [9], the basis for Sec. IV-B). Investigation of such non-convex penalties is left for future research.

APPENDIX A

PROOF OF THEOREM 1

With fixed Γ , let $\bar{G}_1(\Omega)$ denote $\bar{G}(\Omega, \{\Phi\}, \{\Phi^*\})$ up to some irrelevant constants. Then

$$\bar{G}_1(\Omega) = -\frac{1}{p} \ln(|\Omega|) + B \text{tr}(\Sigma^\circ \Omega), \quad (49)$$

where $B = \text{tr}(\bar{S}_k^\circ \Phi_k^*) / (2Mqp)$. We have

$$\mathbf{0} = \frac{\partial \bar{G}_1(\Omega)}{\partial \Omega} = -\frac{1}{p} \Omega^{-1} + B \Sigma^\circ, \quad (50)$$

establishing (33) if $B \neq 0$. The solution is unique since the Hessian of $\bar{G}_1(\Omega)$, given by $\frac{1}{p} \Omega^{-1} \otimes \Omega^{-1}$, is positive definite at $\Omega = \bar{\Omega}(\Gamma)$. Similarly, with fixed Ω , let $\bar{G}_2(\Gamma)$ denote $\bar{G}(\Omega, \{\Phi\}, \{\Phi^*\})$ up to some irrelevant constants. Then

$$\bar{G}_2(\Gamma) = \sum_{k=1}^M \bar{G}_{2k}(\Phi_k), \quad (51)$$

$$\begin{aligned} \bar{G}_{2k}(\Phi_k) &= -\ln(|\Phi_k|) - \ln(|\Phi_k^*|) \\ &+ \frac{1}{p} (\text{tr}(\bar{S}_k^\circ \Phi_k) + \text{tr}(\bar{S}_k^\circ \Phi_k^*)) \text{tr}(\Sigma^\circ \Omega). \end{aligned} \quad (52)$$

The cost $\bar{G}_2(\Gamma)$ is separable in k , Φ_k . We have

$$\mathbf{0} = \frac{\partial \bar{G}_{2k}(\Phi_k)}{\partial \Phi_k^*} = -\Phi_k^{-1} + \frac{1}{p} \bar{S}_k^\circ \text{tr}(\Sigma^\circ \Omega), \quad (53)$$

establishing (34) if $\text{tr}(\Sigma^\circ \Omega) \neq 0$. Similar to [9, Lemma 4], the Hessian of $\bar{G}_{2k}(\Phi_k)$ is positive definite at $\Phi_k = \bar{\Phi}_k(\Omega)$. Therefore, the solution is unique. \square

APPENDIX B

TECHNICAL LEMMAS AND PROOF OF THEOREM 2

Lemma 1 is a restatement of [22, Lemma S.1, Supplementary].

Lemma 1. Assume that i.i.d. data $\mathbf{X}_i \in \mathbb{R}^{p \times q}$, $i = 1, 2, \dots, n$, follows the matrix-valued normal distribution $\mathcal{MVN}(\mathbf{0}, \Sigma^\circ, \Psi^\circ)$, with $\Sigma^\circ \in \mathbb{R}^{p \times p}$, $\Psi^\circ \in \mathbb{R}^{q \times q}$, $\Sigma^\circ \succ \mathbf{0}$ and $\Psi^\circ \succ \mathbf{0}$, i.e., $\text{vec}(\mathbf{X}_i) \sim \mathcal{N}_r(\mathbf{0}, \Psi^\circ \otimes \Sigma^\circ)$. Assume that $\phi_{\max}(\Sigma^\circ) \leq C_{1h} < \infty$ and $\phi_{\max}(\Psi^\circ) \leq C_{2h} < \infty$ for some positive constants C_{1h} and C_{2h} . For any symmetric positive-definite $\Omega \in \mathbb{R}^{p \times p}$ such that $\|\Omega - \Omega^\circ\|_F \leq \gamma$, $\Omega^\circ = (\Sigma^\circ)^{-1}$, we have

$$\begin{aligned} P\left(\max_{i,j} \left| \left[\frac{1}{np} \sum_{i=1}^n \mathbf{X}_i^\top \Omega \mathbf{X}_i - \frac{1}{p} E\{\mathbf{X}_i^\top \Omega \mathbf{X}_i\} \right]_{ij} \right| \geq \delta \right) \\ \leq 4q^2 \left[\exp \left\{ -\frac{np}{2} \left[\frac{\delta}{8(1+\gamma C_{1h})C_{2h}} - \frac{2}{\sqrt{np}} \right]^2 \right\} \right. \\ \left. + \exp \left\{ -\frac{np}{2} \right\} \right] \end{aligned} \quad (54)$$

for any $\delta > 16(1 + \gamma C_{1h})C_{2h}/\sqrt{np}$ •

The lower bound on δ follows from [22, Lemma S.12, Supplementary] and (54) is [22, Eqn. (S.26), Supplementary] in our notation.

Lemma 2 collects some useful results from [42, Theorem 2.3.5].

Lemma 2. Suppose $\mathbf{X} \sim \mathcal{MVN}(\mathbf{0}, \Sigma, \Psi)$ where $\mathbf{X} \in \mathbb{R}^{p \times q}$, $\Sigma \in \mathbb{R}^{p \times p}$, $\Psi \in \mathbb{R}^{q \times q}$, i.e., $\text{vec}(\mathbf{X}) \sim \mathcal{N}_r(\mathbf{0}, \Psi \otimes \Sigma)$. Then

- (i) $\mathbf{X}^\top \sim \mathcal{MVN}(\mathbf{0}, \Psi, \Sigma)$, i.e., $\text{vec}(\mathbf{X}^\top) \sim \mathcal{N}_r(\mathbf{0}, \Sigma \otimes \Psi)$.
- (ii) For any $\mathbf{A} \in \mathbb{R}^{q \times q}$, $E\{\mathbf{X} \mathbf{A} \mathbf{X}^\top\} = \text{tr}(\mathbf{A}^\top \Psi) \Sigma$.
- (iii) For any $\mathbf{B} \in \mathbb{R}^{p \times p}$, $E\{\mathbf{X}^\top \mathbf{B} \mathbf{X}\} = \text{tr}(\mathbf{B}^\top \Sigma) \Psi$.
- (iv) For any $\mathbf{C} \in \mathbb{R}^{q \times p}$, $E\{\mathbf{X} \mathbf{C} \mathbf{X}\} = \Sigma \mathbf{C}^\top \Psi$ •

Lemma 3. Suppose $\mathbf{X} \in \mathbb{C}^{p \times q}$, $\text{vec}(\mathbf{X}) \sim \mathcal{N}_c(\mathbf{0}, \mathbf{S} \otimes \Sigma)$ where $\Sigma \in \mathbb{R}^{p \times p}$, $\mathbf{S} \in \mathbb{C}^{q \times q}$, $\Sigma = \Sigma^\top \succ \mathbf{0}$, $\mathbf{S} = \mathbf{S}^H \succ \mathbf{0}$ and $\mathbf{S} = \mathbf{S}_r + j\mathbf{S}_i$ with $\mathbf{S}_r, \mathbf{S}_i \in \mathbb{R}^{q \times q}$.

- (i) Let $\tilde{\mathbf{X}} = \mathbf{X}_r + j\mathbf{X}_i$, $\mathbf{X}_r, \mathbf{X}_i \in \mathbb{R}^{p \times q}$. Then

$$\tilde{\mathbf{X}} = [\mathbf{X}_r \ \mathbf{X}_i] \sim \mathcal{MVN}(\mathbf{0}, \Sigma, \tilde{\mathbf{S}}) \quad (55)$$

i.e., $\text{vec}(\tilde{\mathbf{X}}) \sim \mathcal{N}_r(\mathbf{0}, \tilde{\mathbf{S}} \otimes \Sigma)$, where

$$\tilde{\mathbf{S}} = \frac{1}{2} \begin{bmatrix} \mathbf{S}_r & -\mathbf{S}_i \\ \mathbf{S}_i & \mathbf{S}_r \end{bmatrix} \in \mathbb{R}^{2q \times 2q}. \quad (56)$$

- (ii) For any $\Omega \in \mathbb{R}^{p \times p}$, $E\{\tilde{\mathbf{X}}^\top \Omega \tilde{\mathbf{X}}\} = \text{tr}(\Omega^\top \Sigma) \tilde{\mathbf{S}}$.
- (iii) For any $\Phi \in \mathbb{C}^{q \times q}$, $\Phi = \Phi_r + j\Phi_i = \Phi^H$, $\Phi_r, \Phi_i \in \mathbb{R}^{q \times q}$,

$$E\{\text{Re}(\mathbf{X} \Phi^* \mathbf{X}^H)\} = E\{\tilde{\mathbf{X}} \tilde{\Phi} \tilde{\mathbf{X}}^\top\} = \text{tr}(\tilde{\Phi}^\top \tilde{\mathbf{S}}) \Sigma \quad (57)$$

where

$$\tilde{\Phi} = \begin{bmatrix} \Phi_r & -\Phi_i \\ \Phi_i & \Phi_r \end{bmatrix} \in \mathbb{R}^{2q \times 2q} \quad \bullet \quad (58)$$

Proof.

- (i) If $\mathbf{x} = \text{vec}(\mathbf{X}) \sim \mathcal{N}_c(\mathbf{0}, \mathbf{S} \otimes \Sigma)$, then by [38, Sec. 2.3],

$$\tilde{\mathbf{x}} = \text{vec}(\tilde{\mathbf{X}}) \sim \mathcal{N}_r(\mathbf{0}, \mathbf{R}) \quad (59)$$

where, with $\mathbf{x} = \mathbf{x}_r + j\mathbf{x}_i$, $\mathbf{x}_r, \mathbf{x}_i \in \mathbb{R}^{pq}$,

$$\mathbf{R} = \begin{bmatrix} E\{\mathbf{x}_r \mathbf{x}_r^\top\} & E\{\mathbf{x}_i \mathbf{x}_r^\top\} \\ E\{\mathbf{x}_r \mathbf{x}_i^\top\} & E\{\mathbf{x}_i \mathbf{x}_i^\top\} \end{bmatrix} = \begin{bmatrix} \mathbf{R}_{rr} & \mathbf{R}_{ir} \\ \mathbf{R}_{ri} & \mathbf{R}_{ii} \end{bmatrix} \quad (60)$$

$$\mathbf{R}_{rr} = \mathbf{R}_{ii}, \quad \mathbf{R}_{ri} = -\mathbf{R}_{ri}^\top = \mathbf{R}_{ir}^\top. \quad (61)$$

Now $\mathbf{R}_{rr} = \frac{1}{2}\mathbf{S}_r \otimes \Sigma = \mathbf{R}_{ii}$ and $\mathbf{R}_{ri} = -\frac{1}{2}\mathbf{S}_i^\top \otimes \Sigma$. Therefore, $\mathbf{R} = \tilde{\mathbf{S}} \otimes \Sigma$, yielding the desired result.

- (ii) It follows from Lemma 2(iii) and Lemma 3(i).
- (iii) Since $\Phi = \Phi^H$, it follows that $\Phi_r = \Phi_r^\top$ and $\Phi_i = -\Phi_i^\top$. We have $\text{Re}(\mathbf{X}\Phi^*\mathbf{X}^H) = \tilde{\mathbf{X}}\tilde{\Phi}\tilde{\mathbf{X}}^\top$. Then the given expression for $E\{\tilde{\mathbf{X}}\tilde{\Phi}\tilde{\mathbf{X}}^\top\}$ follows from Lemma 2(ii). \square

We now consider a tail bound on $\tilde{\Theta}_k$ defined in (27). First we need Lemma 4.

Lemma 4. Given $\mathbf{S} \in \mathbb{C}^{q \times q}$ and $\tilde{\mathbf{S}} \in \mathbb{R}^{2q \times 2q}$ as in Lemma 3. Then $\tilde{\mathbf{S}} \succ \mathbf{0}$ and $\phi_{\max}(\tilde{\mathbf{S}}) = \frac{1}{2}\phi_{\max}(\mathbf{S})$.

Proof. If λ is an eigenvalue of \mathbf{S} , then for some $\mathbf{v} = \mathbf{v}_r + j\mathbf{v}_i \in \mathbb{C}^q$, $\mathbf{v}_r, \mathbf{v}_i \in \mathbb{R}^q$, we have $\mathbf{S}\mathbf{v} = \lambda\mathbf{v}$, where λ is real positive since \mathbf{S} is Hermitian, positive-definite. It then follows that

$$\tilde{\mathbf{S}} \begin{bmatrix} \mathbf{v}_r \\ \mathbf{v}_i \end{bmatrix} = \frac{1}{2}\lambda \begin{bmatrix} \mathbf{v}_r \\ \mathbf{v}_i \end{bmatrix}, \quad \tilde{\mathbf{S}} \begin{bmatrix} -\mathbf{v}_i \\ \mathbf{v}_r \end{bmatrix} = \frac{1}{2}\lambda \begin{bmatrix} -\mathbf{v}_i \\ \mathbf{v}_r \end{bmatrix}. \quad (62)$$

That is, each eigenvalue of \mathbf{S} is also an eigenvalue of $2\tilde{\mathbf{S}}$ with multiplicity two. This proves the desired result. \square

Lemma 5. Under assumptions (A1) and (A2), for any symmetric positive-definite $\hat{\Omega} \in \mathbb{R}^{p \times p}$ such that $\|\hat{\Omega} - \Omega^\circ\|_F \leq \gamma_p$, $\Omega^\circ = (\Sigma^\circ)^{-1}$, and $\tau > 2$, we have

$$\begin{aligned} P\left(\max_{k,i,j} |[\tilde{\Theta}_k^* - E\{\tilde{\Theta}_k^*\}]_{ij}| \geq C_{0q} \sqrt{\frac{\ln(M^{1/\tau}q)}{Kp}}\right) \\ \leq \frac{1}{q^{\tau-2}} + 16Mq^2 e^{-Kp/2} \end{aligned} \quad (63)$$

for any $q \geq 1$, where C_{0q} is given by (37) and

$$E\{\tilde{\Theta}_k^*\} = \frac{1}{p} \text{tr}(\hat{\Omega}\Sigma^\circ) (\tilde{\mathbf{S}}_k^\circ)^*. \quad (64)$$

Proof. Let $\mathbf{D}_z(\tilde{f}_{k,\ell}) = \mathbf{D}_{r,kl} + j\mathbf{D}_{i,kl}$, $\mathbf{D}_{r,kl}, \mathbf{D}_{i,kl} \in \mathbb{R}^{p \times q}$. Define

$$\mathbf{X}_{kl} = [\mathbf{D}_{r,kl} \quad \mathbf{D}_{i,kl}] \in \mathbb{R}^{p \times (2q)}, \quad (65)$$

$$\mathbf{B}_{kl} = \mathbf{X}_{kl}^\top \hat{\Omega} \mathbf{X}_{kl}, \quad \mathbf{F}_k = \frac{1}{Kp} \sum_{\ell=-m_t}^{m_t} \mathbf{B}_{kl}. \quad (66)$$

Since

$$\begin{aligned} \mathbf{D}_z^H(\tilde{f}_{k,\ell}) \hat{\Omega} \mathbf{D}_z(\tilde{f}_{k,\ell}) &= \mathbf{D}_{r,kl}^\top \hat{\Omega} \mathbf{D}_{r,kl} + \mathbf{D}_{i,kl}^\top \hat{\Omega} \mathbf{D}_{i,kl} \\ &\quad + j[\mathbf{D}_{r,kl}^\top \hat{\Omega} \mathbf{D}_{i,kl} - \mathbf{D}_{i,kl}^\top \hat{\Omega} \mathbf{D}_{r,kl}], \end{aligned} \quad (67)$$

it follows that

$$\max_{k,i,j} |[\tilde{\Theta}_k^* - E\{\tilde{\Theta}_k^*\}]_{ij}| \leq 4 \max_{k,i,j} |[\mathbf{F}_k - E\{\mathbf{F}_k\}]_{ij}|. \quad (68)$$

Therefore,

$$\begin{aligned} \left\{ \max_{k,i,j} |[\mathbf{F}_k - E\{\mathbf{F}_k\}]_{ij}| < \frac{\delta}{4} \right\} \\ \subseteq \left\{ \max_{k,i,j} |[\tilde{\Theta}_k^* - E\{\tilde{\Theta}_k^*\}]_{ij}| < \delta \right\}, \end{aligned} \quad (69)$$

implying

$$\begin{aligned} P\left(\max_{k,i,j} |[\tilde{\Theta}_k^* - E\{\tilde{\Theta}_k^*\}]_{ij}| \geq \delta\right) \\ \leq P\left(\max_{k,i,j} |[\mathbf{F}_k - E\{\mathbf{F}_k\}]_{ij}| \geq \frac{\delta}{4}\right). \end{aligned} \quad (70)$$

Since $\mathbf{d}_z(\tilde{f}_{k,\ell}) = \text{vec}(\mathbf{D}_z(\tilde{f}_{k,\ell})) \sim \mathcal{N}_c(\mathbf{0}, \bar{\mathbf{S}}^\circ(\tilde{f}_k) \otimes \Sigma^\circ)$, it follows from Lemma 3(i) that

$$\mathbf{X}_{kl} \sim \mathcal{M}\mathcal{V}\mathcal{N}(\mathbf{0}, \Sigma^\circ, \tilde{\mathbf{S}}_k^\circ), \quad (71)$$

$$\tilde{\mathbf{S}}_k^\circ = \frac{1}{2} \begin{bmatrix} \bar{\mathbf{S}}_{rk}^\circ & -\bar{\mathbf{S}}_{ik}^\circ \\ \bar{\mathbf{S}}_{ik}^\circ & \bar{\mathbf{S}}_{rk}^\circ \end{bmatrix}, \quad \bar{\mathbf{S}}^\circ(\tilde{f}_k) = \bar{\mathbf{S}}_{rk}^\circ + j\bar{\mathbf{S}}_{ik}^\circ. \quad (72)$$

By assumption (A4), $\phi_{\max}(\Sigma^\circ) \leq \beta_{p,\max}$ and additionally, by Lemma 4, $\phi_{\max}(\tilde{\mathbf{S}}_k^\circ) \leq \beta_{q,\max}/2$ for every k . With $a = 4(1 + \gamma_p\beta_{p,\max})\beta_{q,\max}$, invoking Lemma 1 for the sum \mathbf{F}_k , we have

$$\begin{aligned} P\left(\max_{i,j} |[\mathbf{F}_k - E\{\mathbf{F}_k\}]_{ij}| \geq \frac{\delta}{4}\right) \\ \leq 4(2q)^2 \left[\exp\left\{-\frac{Kp}{2} \left[\frac{\delta/4}{a} - \frac{2}{\sqrt{Kp}}\right]^2\right\} + e^{-Kp/2} \right] = P_{qtb}. \end{aligned} \quad (73)$$

Maximizing over all $k = 1, 2, \dots, M$, and using the union bound, we obtain

$$P\left(\max_{k,i,j} |[\mathbf{F}_k - E\{\mathbf{F}_k\}]_{ij}| \geq \frac{\delta}{4}\right) \leq MP_{qtb}. \quad (74)$$

For $\tau > 2$, pick $\delta = 4a(\sqrt{2 \ln(16Mq^\tau)/(Kp)} + 2/\sqrt{Kp})$, leading to $\delta = C_{0q}\sqrt{\ln(M^{1/\tau}q)/(Kp)}$ and $(Kp/2)[(\delta/(4a)) - 2/\sqrt{Kp}]^2 = \ln(16Mq^\tau)$. Thus

$$\begin{aligned} MP_{qtb} &= 16Mq^2 \left[e^{-\ln(16Mq^\tau)} + e^{-Kp/2} \right] \\ &= \frac{1}{q^{\tau-2}} + 16Mq^2 e^{-Kp/2}. \end{aligned} \quad (75)$$

Thus we have established (63). The lower bound on $\delta/4$ specified in Lemma 1 is satisfied if $(\delta/(4a)) > 2/\sqrt{Kp}$, which is true for any $q \geq 1$. Turning to (64), by (66), (71) and Lemma 2(iii), we have

$$E\{\mathbf{B}_{kl}\} = \text{tr}(\hat{\Omega}\Sigma^\circ) \frac{1}{2} \begin{bmatrix} \bar{\mathbf{S}}_{rk}^\circ & -\bar{\mathbf{S}}_{ik}^\circ \\ \bar{\mathbf{S}}_{ik}^\circ & \bar{\mathbf{S}}_{rk}^\circ \end{bmatrix}. \quad (76)$$

By assumption (A1), (66), (67) and (76),

$$\begin{aligned} E\{\mathbf{D}_z^H(\tilde{f}_{k,\ell}) \hat{\Omega} \mathbf{D}_z(\tilde{f}_{k,\ell})\} &= \frac{1}{2} \text{tr}(\hat{\Omega}\Sigma^\circ) (2\bar{\mathbf{S}}_{rk}^\circ - j2\bar{\mathbf{S}}_{ik}^\circ) \\ &= \text{tr}(\hat{\Omega}\Sigma^\circ) (\tilde{\mathbf{S}}_k^\circ)^*. \end{aligned} \quad (77)$$

By (27) and (77), we obtain (64). \square

Now we consider a tail bound on $\tilde{\Theta}$ defined in (25).

Lemma 6. Under assumptions (A1) and (A2), for any Hermitian positive-definite $\hat{\Phi}_k \in \mathbb{C}^{q \times q}$, $k = 1, 2, \dots, M$, such that $\|\hat{\Gamma} - \Gamma^\circ\|_F \leq \gamma_q$, $\Gamma^\circ = [\Phi_1^\circ, \dots, \Phi_M^\circ]$, $\hat{\Gamma} = [\hat{\Phi}_1, \dots, \hat{\Phi}_M]$, and $\tau > 2$, we have

$$\begin{aligned} P\left(\max_{i,j} |[\tilde{\Theta} - E\{\tilde{\Theta}\}]_{ij}| \geq C_{0p} \sqrt{\frac{\ln(p)}{KqM}}\right) \\ \leq \frac{1}{p^{\tau-2}} + 4p^2 e^{-KqM} \end{aligned} \quad (78)$$

for any $p \geq 1$ where, where C_{0p} is given by (40), and

$$E\{\check{\Theta}\} = \left[\frac{1}{2Mq} \sum_{k=1}^M \text{tr}(\bar{\mathbf{S}}_k^\circ \hat{\Phi}_k + (\bar{\mathbf{S}}_k^\circ \hat{\Phi}_k)^*) \right] \Sigma^\circ. \quad (79)$$

Proof. We have

$$\text{Re}(D_z(\tilde{f}_{k,\ell}) \hat{\Phi}_k D_z^H(\tilde{f}_{k,\ell})) = \mathbf{X}_{kl} \tilde{\Phi}_k \mathbf{X}_{kl}^\top, \quad (80)$$

where \mathbf{X}_{kl} is as in (65) and

$$\tilde{\Phi}_k = \begin{bmatrix} \hat{\Phi}_{rk} & -\hat{\Phi}_{ik} \\ \hat{\Phi}_{ik} & \hat{\Phi}_{rk} \end{bmatrix} \in \mathbb{R}^{2q \times 2q}. \quad (81)$$

Define

$$\check{\Phi} = \begin{bmatrix} \tilde{\Phi}_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \tilde{\Phi}_2 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \tilde{\Phi}_M \end{bmatrix} \in \mathbb{R}^{(2qM) \times (2qM)}, \quad (82)$$

$$\check{\mathbf{X}}_\ell = [\mathbf{X}_{1\ell} \quad \mathbf{X}_{2\ell} \quad \cdots \quad \mathbf{X}_{M\ell}]^\top. \quad (83)$$

Then we can express $\check{\Theta}$ as

$$\check{\Theta} = \frac{1}{MKq} \sum_{\ell=-m_t}^{m_t} \check{\mathbf{X}}_\ell^\top \check{\Phi} \check{\mathbf{X}}_\ell. \quad (84)$$

Since $\mathbf{X}_{kl} \sim \mathcal{M}\mathcal{V}\mathcal{N}(\mathbf{0}, \Sigma^\circ, \tilde{\mathbf{S}}_k^\circ)$, and \mathbf{X}_{k_1l} and \mathbf{X}_{k_2l} are independent for $k_1 \neq k_2$, we have

$$\check{\mathbf{X}}_\ell^\top \sim \mathcal{M}\mathcal{V}\mathcal{N}(\mathbf{0}, \Sigma^\circ, \check{\mathbf{S}}^\circ), \quad \check{\mathbf{X}}_\ell \sim \mathcal{M}\mathcal{V}\mathcal{N}(\mathbf{0}, \check{\mathbf{S}}^\circ, \Sigma^\circ), \quad (85)$$

where

$$\check{\mathbf{S}}^\circ = \begin{bmatrix} \tilde{\mathbf{S}}_1^\circ & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \tilde{\mathbf{S}}_2^\circ & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \tilde{\mathbf{S}}_M^\circ \end{bmatrix} \in \mathbb{R}^{(2qM) \times (2qM)}. \quad (86)$$

By assumption (A4), $\phi_{\max}(\Sigma^\circ) \leq \beta_{p,\max}$ and additionally, by Lemma 4, $\phi_{\max}(\check{\mathbf{S}}^\circ) \leq \beta_{q,\max}/2$. With $b = 8(1 + \gamma_q \beta_{q,\max}/2) \beta_{p,\max}$, apply Lemma 1 to the sum $\frac{1}{2} \check{\Theta}$ to obtain

$$\begin{aligned} & P\left(\max_{i,j} \left| \frac{1}{2} [\check{\Theta} - E\{\check{\Theta}\}]_{ij} \right| \geq \delta\right) \\ & \leq 4p^2 \left[\exp\left\{-\frac{2qMK}{2} \left[\frac{\delta}{b} - \frac{2}{\sqrt{2qMK}} \right]^2\right\} + e^{-2qMK/2} \right] \\ & = P_{ptb}. \end{aligned} \quad (87)$$

For $\tau > 2$, pick $\delta = b(\sqrt{\ln(4p^\tau)/(KqM)} + \sqrt{2/(KqM)})$, leading to $(2qMK/2)[(\delta/b) - \sqrt{2/(KqM)}]^2 = \ln(4p^\tau)$. Thus

$$P_{ptb} = 4p^2 \left[e^{-\ln(4p^\tau)} + e^{-qMK} \right] = \frac{1}{p^{\tau-2}} + 4p^2 e^{-qMK}. \quad (88)$$

The lower bound on δ specified in Lemma 1 is satisfied if $(\delta/b) > \sqrt{2/(KqM)}$, which is true for any $p \geq 1$. With our choice of δ , we have $2\delta = C_{0p} \sqrt{\frac{\ln(p)}{KqM}}$, establishing (78). Turning to (79), by (85) and Lemma 2(iii), we have

$$E\{\check{\mathbf{X}}_\ell^\top \check{\Phi} \check{\mathbf{X}}_\ell\} = \text{tr}(\check{\Phi}^\top \check{\mathbf{S}}^\circ) \Sigma^\circ = \text{tr}\left(\sum_{k=1}^M \check{\Phi}_k^\top \tilde{\mathbf{S}}_k^\circ\right) \Sigma^\circ. \quad (89)$$

By (72) and (81),

$$\begin{aligned} \text{tr}(\check{\Phi}_k^\top \tilde{\mathbf{S}}_k^\circ) &= \text{tr}(\hat{\Phi}_{rk} \bar{\mathbf{S}}_{rk}^\circ + \hat{\Phi}_{ik} \bar{\mathbf{S}}_{ik}^\circ) \\ &= \frac{1}{2} \text{tr}(\bar{\mathbf{S}}_k^\circ \hat{\Phi}_k + (\bar{\mathbf{S}}_k^\circ \hat{\Phi}_k)^*). \end{aligned} \quad (90)$$

Using (84), (89) and (90), we have (79). \square

Proof of Theorem 2(i). Let $\Omega = \bar{\Omega}(\Gamma) + \Delta$ with $\Omega, \bar{\Omega}(\Gamma) \succ \mathbf{0}$, and denote $Q(\Omega) = L_1(\Omega) - L_1(\bar{\Omega}(\Gamma))$. For the rest of the proof, we will denote $\bar{\Omega}(\Gamma)$ by $\bar{\Omega}$. Then $\hat{\Omega}(\Gamma)$ minimizes $Q(\Omega)$, or equivalently, $\hat{\Delta} = \hat{\Omega}(\Gamma) - \bar{\Omega}$ minimizes $J(\Delta) = Q(\bar{\Omega} + \Delta)$. Consider the set

$$\Psi_p(R_p) := \{\Delta : \Delta = \Delta^\top, \|\Delta\|_F = R_p r_{pn}\} \quad (91)$$

where $R_p = 17C_{0p}/\beta_{p,\min}^2$ and r_{pn} is as in (42). Since $J(\hat{\Delta}) \leq J(\mathbf{0}) = 0$, if we can show that $\inf_{\Delta} \{J(\Delta) : \Delta \in \Psi_p(R_p)\} > 0$, then the minimizer $\hat{\Delta}$ must be inside the sphere defined by $\Psi_p(R_p)$, and hence, $\|\hat{\Delta}\|_F \leq R_p r_{pn}$. It is shown in [41, (9)] that $\ln(|\bar{\Omega} + \Delta|) - \ln(|\bar{\Omega}|) = \text{tr}(\bar{\Omega}^{-1} \Delta) - \tilde{B}_1$ where, with $H(\bar{\Omega}, \Delta, v) = (\bar{\Omega} + v\Delta)^{-1} \otimes (\bar{\Omega} + v\Delta)^{-1}$ and v denoting a real scalar,

$$\tilde{B}_1 := \text{vec}(\Delta)^\top \left(\int_0^1 (1-v) H(\bar{\Omega}, \Delta, v) dv \right) \text{vec}(\Delta). \quad (92)$$

We have

$$J(\Delta) = \sum_{i=1}^3 B_i, \quad B_1 = \frac{1}{p} \tilde{B}_1, \quad (93)$$

$$B_2 := \frac{1}{p} \text{tr}((\check{\Theta} - \bar{\Omega}^{-1}) \Delta), \quad (94)$$

$$B_3 := \lambda_p(\|\bar{\Omega}^- + \Delta^- \|_1 - \|\bar{\Omega}^- \|_1). \quad (95)$$

By (33) and (79), $\bar{\Omega}^{-1} = E\{\check{\Theta}\} = \left(\frac{1}{2Mq} \sum_{k=1}^M \text{tr}(\bar{\mathbf{S}}_k^\circ \hat{\Phi}_k + (\bar{\mathbf{S}}_k^\circ \hat{\Phi}_k)^*) \right) \Sigma^\circ$ (where we replaced $\hat{\Phi}_k$ with Φ_k). By Lemma 6, $\max_{i,j} |[\check{\Theta} - E\{\check{\Theta}\}]_{ij}| \geq C_{0p} \sqrt{\frac{\ln(p)}{KqM}}$ w.h.p. (which refers to $1 - \frac{1}{p^{\tau-2}} - 4p^2 e^{-KqM}$, cf. (78)). Following [41, p. 502], we have

$$\tilde{B}_1 \geq \|\Delta\|_F^2 / (2(\|\bar{\Omega}\| + \|\Delta\|)^2). \quad (96)$$

Turning to $E\{\check{\Theta}\}$, we have

$$\text{tr}(\bar{\mathbf{S}}_k^\circ \Phi_k + (\bar{\mathbf{S}}_k^\circ \Phi_k)^*) = 2\text{Re} \text{tr}(\bar{\mathbf{S}}_k^\circ (\Phi_k - \Phi_k^\circ + \Phi_k^\circ)) \quad (97)$$

$$= 2\text{Re} \text{tr}(\bar{\mathbf{S}}_k^\circ (\Phi_k - \Phi_k^\circ)) + 2\text{tr}(\mathbf{I}_q) \quad (98)$$

$$\geq 2q - 2 \|\bar{\mathbf{S}}_k^\circ\|_F \|\Phi_k - \Phi_k^\circ\|_F \quad (99)$$

where we used $|\text{tr}(BC^H)| \leq \|B\|_F \|C\|_F$ (Cauchy-Schwarz inequality). Since $\|\bar{\mathbf{S}}_k^\circ\|_F \leq \sqrt{q} \|\tilde{\mathbf{S}}_k^\circ\| \leq \sqrt{q} \beta_{q,\max}$ and $\sum_{k=1}^M \|\Phi_k - \Phi_k^\circ\|_F \leq \sqrt{M} \|\Gamma - \Gamma^\circ\|_F \leq \sqrt{M} \gamma_q$, we have

$$A = 2\text{Re} \sum_{k=1}^M \text{tr}(\bar{\mathbf{S}}_k^\circ \Phi_k) \geq 2Mq - 2\sqrt{Mq} \beta_{q,\max} \gamma_q \quad (100)$$

$$\geq 2Mq - 2Mq \beta_{q,\max} \gamma_q = 1.8Mq, \quad (101)$$

where we have used the facts that $\sqrt{Mq} \leq Mq$ and $\beta_{q,\max} \gamma_q = 0.1$, as defined in (38). Therefore, $\|\bar{\Omega}^{-1}\| = \|E\{\check{\Theta}\}\| \geq 0.9 \|\Sigma^\circ\|$, implying $\|\bar{\Omega}\| \leq 10/(9 \|\Sigma^\circ\|) \leq 10/(9 \beta_{p,\min}) \leq 1.5/\beta_{p,\min}$. Using (93), (96), and the facts

$\|\bar{\Omega}\| \leq 1.5/\beta_{p,\min}$ and $\|\Delta\| \leq \|\Delta\|_F = R_p r_{pn}$, we obtain w.h.p.

$$B_1 \geq \|\Delta\|_F^2 \beta_{p,\min}^2 / (8p), \quad (102)$$

for $n > N_p$, since $r_{pn} \leq \beta_{p,\min}/(34C_{0p})$ for $n > N_p$ and $R_p r_{pn} \leq 0.5/\beta_{p,\min}$.

We now consider B_2 given by (94). Define $\bar{S}_p = S_p \cup \{\{i, j\} : i = j\}$ so that $|\bar{S}_p| = s_p + p$. We have

$$|B_2| \leq B_{12} + B_{22}, \quad pB_{12} = \left| \sum_{\{i,j\} \in \bar{S}_p} [\check{\Theta} - \bar{\Omega}^{-1}]_{ij} \Delta_{ji} \right|,$$

$$pB_{22} = \left| \sum_{\{i,j\} \in \bar{S}_p^c} [\check{\Theta} - \bar{\Omega}^{-1}]_{ij} \Delta_{ji} \right|,$$

where \bar{S}_p^c denotes the complement of set \bar{S}_p . For an index set B and a matrix $C \in \mathbb{R}^{p \times p}$, we write C_B to denote a matrix in $\mathbb{R}^{p \times p}$ such that $[C_B]_{ij} = C_{ij}$ if $(i, j) \in B$, and $[C_B]_{ij} = 0$ if $(i, j) \notin B$. Using $|\sum_{\{i,j\} \in \bar{S}_p} \Delta_{ij}| \leq \sqrt{s_p + p} \|\Delta\|_F$ (by Cauchy-Schwarz inequality),

$$pB_{12} \leq \max_{i,j} |[\check{\Theta} - \bar{\Omega}^{-1}]_{ij}| \sum_{\{i,j\} \in \bar{S}_p} |\Delta_{ij}|,$$

$$\leq C_{0p} \sqrt{\ln(p)/(KqM)} \sqrt{s_p + p} \|\Delta\|_F = C_{0p} r_{pn} \|\Delta\|_F. \quad (103)$$

We will combine B_{22} with B_3 . By (95),

$$B_3 = \lambda_p (\|\bar{\Omega}^- + \Delta_{S_p}^- \|_1 + \|\Delta_{\bar{S}_p^c}^- \|_1 - \|\bar{\Omega}^- \|_1)$$

$$\geq \lambda_p (\|\Delta_{\bar{S}_p^c}^- \|_1 - \|\Delta_{S_p}^- \|_1), \quad (104)$$

using the triangle inequality $\|\bar{\Omega}^- + \Delta_{S_p}^- \|_1 \geq \|\bar{\Omega}^- \|_1 - \|\Delta_{S_p}^- \|_1$ and the fact $\bar{\Omega}_{\bar{S}_p^c}^- = \bar{\Omega}_{\bar{S}_p^c}^- = \mathbf{0}$. Hence, $B_2 + B_3 \geq -B_{12} - B_{22} + \lambda_p (\|\Delta_{\bar{S}_p^c}^- \|_1 - \|\Delta_{S_p}^- \|_1)$. But $pB_{22} \leq C_{0p} \sqrt{\ln(p)/(KqM)} \|\Delta_{\bar{S}_p^c}^- \|_1$ w.h.p., therefore,

$$B_2 + B_3 \geq (\lambda_p - C_{0p} \sqrt{\ln(p)/(p^2 KqM)}) \|\Delta_{\bar{S}_p^c}^- \|_1$$

$$- \lambda_p \|\Delta_{S_p}^- \|_1 - C_{0p} r_{pn} \|\Delta\|_F / p. \quad (105)$$

Using the fact that by (43), the first term on right side of (105) is nonnegative, and $\|\Delta_{S_p}^- \|_1 \leq \sqrt{s_p} \|\Delta\|_F$ by the Cauchy-Schwarz inequality, we obtain $B_2 + B_3 \geq -(\lambda_p \sqrt{s_p} + r_{pn}/p) \|\Delta\|_F$. Thus, by (43), (93) and (102)

$$J(\Delta) \geq \frac{\|\Delta\|_F^2 \beta_{p,\min}^2}{8p} - (\lambda_p \sqrt{s_p} + C_{0p} r_{pn}/p) \|\Delta\|_F$$

$$\geq \frac{\|\Delta\|_F^2 \beta_{p,\min}^2}{8p} - \frac{2C_{0p} r_{pn} \|\Delta\|_F}{p}$$

$$= \frac{\|\Delta\|_F^2 \beta_{p,\min}^2}{8p} \left(1 - \frac{16}{17}\right) > 0 \quad (106)$$

using $\|\Delta\|_F = R_p r_{pn}$ and $R_p = 17C_{0p}/\beta_{p,\min}^2$. This proves Theorem 2(i). \square

Proof of Theorem 2(ii). With Γ as in (17), let $\Gamma = \bar{\Gamma}(\Omega) + \Lambda$ with $\Phi_k = \bar{\Phi}_k^H \succ \mathbf{0}$, and denote $Q(\Gamma) = L_2(\Gamma) - L_2(\bar{\Gamma}(\Omega))$. For the rest of the proof, we will denote $\bar{\Gamma}(\Omega)$ by $\bar{\Gamma}$. Then $\hat{\Gamma}(\Omega)$ minimizes $Q(\Gamma)$, or equivalently, $\hat{\Lambda} = \hat{\Gamma}(\Omega) - \bar{\Gamma}$ minimizes $J(\Lambda) = Q(\hat{\Gamma} + \Lambda)$. Note that $\Lambda = [\Lambda_1, \dots, \Lambda_M] \in$

$\mathbb{C}^{q \times (qM)}$ and $\Lambda_k = \Phi_k - \bar{\Phi}_k$, $k = 1, \dots, M$, where $\bar{\Phi}_k = \bar{\Phi}_k(\Omega) = p\Phi_k^{\circ}/\text{tr}(\Sigma^{\circ}\Omega)$ by (34). Consider the set

$$\Psi_q(R_q) := \{\Lambda : \Lambda_k = \Lambda_k^H, k = 1, \dots, M, \|\Lambda\|_F = R_q r_{qn}\} \quad (107)$$

where $R_q = 17C_{0q}/\beta_{q,\min}^2$ and r_{qn} is as in (41). Similar to the proof of Theorem 2(i), our objective is to show that $\inf_{\Lambda} \{J(\Lambda) : \Lambda \in \Psi_q(R_q)\} > 0$, which would ensure $\|\hat{\Lambda}\|_F \leq R_q r_{qn}$ w.h.p. It is shown in [9, Lemma 5] that $\ln(|\bar{\Phi}_k + \Lambda_k|) - \ln(|\bar{\Phi}_k|) + \ln(|\bar{\Phi}_k^* + \Lambda_k^*|) - \ln(|\bar{\Phi}_k^*|) = \text{tr}(\bar{\Phi}_k^{-1} \Lambda_k + (\bar{\Phi}_k^{-1} \Lambda_k)^*) - \bar{B}_{1k}$ where

$$\bar{B}_{1k} = g^H(\Lambda_k) \left(\int_0^1 (1-v) \mathbf{H}_k(\bar{\Phi}_k, \Lambda_k, v) dv \right) g(\Lambda_k), \quad (108)$$

$$g(\Lambda_k) = \begin{bmatrix} \text{vec}(\Lambda_k) \\ \text{vec}(\Lambda_k^*) \end{bmatrix}, \quad \mathbf{H}_k(\bar{\Phi}_k, \Lambda_k, v) = \begin{bmatrix} \mathbf{H}_{11k} & \mathbf{0} \\ \mathbf{0} & \mathbf{H}_{22k} \end{bmatrix}, \quad (109)$$

$$\mathbf{H}_{11k} = (\bar{\Phi}_k + v\Lambda_k)^{-*} \otimes (\bar{\Phi}_k + v\Lambda_k)^{-1}, \quad (110)$$

$$\mathbf{H}_{22k} = (\bar{\Phi}_k + v\Lambda_k)^{-1} \otimes (\bar{\Phi}_k + v\Lambda_k)^{-*}, \quad (111)$$

and v is a real scalar. Therefore,

$$J(\Lambda) = \sum_{k=1}^M \sum_{i=1}^3 B_{ik} + B_4, \quad B_{ik} = \frac{1}{2Mq} \bar{B}_{1k}, \quad (112)$$

$$B_{2k} = \frac{1}{2Mq} \text{tr}(\bar{B}_{2k} + \bar{B}_{2k}^*), \quad \bar{B}_{2k} = (\check{\Theta} - \bar{\Phi}_k^{-1}) \Lambda_k, \quad (113)$$

$$B_{3k} = \alpha \lambda_q (\|\bar{\Phi}_k^- + \Lambda_k^- \|_1 - \|\bar{\Phi}_k^- \|_1), \quad (114)$$

$$B_4 = (1 - \alpha) \sqrt{M} \lambda_q \sum_{i \neq j}^p (\|\Phi^{(ij)} + \Lambda^{(ij)}\| - \|\Phi^{(ij)}\|). \quad (115)$$

By [9, Eqn. (B.43)], we have

$$B_{1k} \geq \frac{1}{2Mq} \frac{\|\Lambda_k\|_F^2}{(\|\bar{\Phi}_k\| + \|\Lambda_k\|)^2}. \quad (116)$$

Now $\text{tr}(\Sigma^{\circ}\Omega) = \text{tr}(\Sigma^{\circ}(\Omega - \Omega^{\circ} + \Omega^{\circ})) = \text{tr}(\Sigma^{\circ}(\Omega - \Omega^{\circ}) + p)$. Since $|\text{tr}(\Sigma^{\circ}(\Omega - \Omega^{\circ}))| \leq \|\Sigma^{\circ}\|_F \|\Omega - \Omega^{\circ}\|_F \leq \sqrt{p} \beta_{p,\max} \gamma_p$, we have $|\text{tr}(\Sigma^{\circ}\Omega)| \geq p - \sqrt{p} \beta_{p,\max} \gamma_p \geq p - p \beta_{p,\min} \gamma_p = 0.9p$ since $\gamma_p = 0.1/\beta_{p,\min}$. Therefore, $\|\bar{\Phi}_k\| \leq p \|\Phi_k^{\circ}\| / (0.9p) \leq 1.5/\beta_{q,\min}$. Also, $\|\Lambda_k\| \leq \|\Lambda_k\|_F \leq \|\Lambda\|_F = R_q r_{qn}$. Therefore,

$$\sum_{k=1}^M B_{1k} \geq \frac{1}{2Mq} \frac{\sum_{k=1}^M \|\Lambda_k\|_F^2}{(1.5/\beta_{q,\min} + R_q r_{qn})^2}$$

$$\geq \frac{\|\Lambda\|_F^2 \beta_{q,\min}^2}{8Mq} \quad (117)$$

w.h.p. for $n > N_q$, since $r_{qn} \leq \beta_{q,\min}/(34C_{0q})$ for $n > N_q$ and $R_q r_{qn} \leq 0.5/\beta_{q,\min}$.

We now bound B_{2k} noting that $|B_{2k}| \leq L_{1k} + L_{2k}$ where

$$L_{1k} = \frac{2}{2Mq} \left| \sum_{\{i,j\} \in \bar{S}_q} [\check{\Theta} - \bar{\Phi}_k^{-1}]_{ij} [\Lambda_k]_{ji} \right|,$$

$$L_{2k} = \frac{2}{2Mq} \left| \sum_{\{i,j\} \in \bar{S}_q^c} [\check{\Theta} - \bar{\Phi}_k^{-1}]_{ij} [\Lambda_k]_{ji} \right|$$

where $\bar{S}_q = \mathcal{S}_q \cup \{\{i, j\} : i = j\}$ so that $|\bar{S}_q| = s_q + q$. Using Lemma 5 and $|\sum_{\{i,j\} \in \bar{S}_q} [\mathbf{A}_k]_{ij}| \leq \sqrt{s_q + q} \|\mathbf{A}_k\|_F$ (by Cauchy-Schwarz inequality), we have

$$\begin{aligned} L_{1k} &\leq \frac{1}{Mq} C_{0q} \sqrt{\frac{\ln(M^{1/\tau} q)}{Kp}} \left| \sum_{\{i,j\} \in \bar{S}_q} [\mathbf{A}_k]_{ij} \right| \\ &\leq \frac{C_{0q}}{M^{3/2} q} r_{qn} \|\mathbf{A}_k\|_F, \end{aligned} \quad (118)$$

$$L_{2k} \leq \frac{C_{0q}}{Mq} \sqrt{\frac{\ln(M^{1/\tau} q)}{Kp}} \|\mathbf{A}_{k\mathcal{S}_q^c}\|_1. \quad (119)$$

Alternatively, as in [9, Eqn. (B.56)], with $B_2 = \sum_{k=1}^M B_{2k}$,

$$\begin{aligned} |B_2| &\leq \frac{2}{2Mq} \sum_{i,j=1}^p \sum_{k=1}^M |[\tilde{\Theta} - \tilde{\Phi}_k^{-1}]_{ij}| |[\mathbf{A}_k]_{ji}| \\ &\leq \frac{C_{0q}}{Mq} \sqrt{\frac{\ln(M^{1/\tau} q)}{Kp}} \sum_{i,j=1}^p \sum_{k=1}^M |[\mathbf{A}_k]_{ij}|. \end{aligned} \quad (120)$$

Define $\check{\mathbf{A}} \in \mathbb{R}^{q \times q}$ with $[\check{\mathbf{A}}]_{ij} = \|\mathbf{A}^{(ij)}\|_F$ and as in (23), $\mathbf{\Lambda}^{(ij)} := [[\mathbf{A}_1]_{ij} \cdots [\mathbf{A}_M]_{ij}]^\top \in \mathbb{C}^M$. Using $\sum_{k=1}^M |[\mathbf{A}_k]_{ij}| \leq \sqrt{M} \|\mathbf{A}^{(ij)}\|_F$, we have

$$|B_2| \leq \frac{C_{0q}}{\sqrt{M} q} \sqrt{\frac{\ln(M^{1/\tau} q)}{Kp}} \|\check{\mathbf{A}}\|_1. \quad (121)$$

Mimicking [9, Eqns. (B.56)-(B.58)], we have $B_{3k} \geq \alpha \lambda_q (\|\mathbf{A}_{k\mathcal{S}_q^c}^-\|_1 - \|\mathbf{A}_{k\mathcal{S}_q}^-\|_1)$ and $B_4 \geq (1 - \alpha) \sqrt{M} \lambda_q (\|\check{\mathbf{A}}_{k\mathcal{S}_q^c}^-\|_1 - \|\check{\mathbf{A}}_{k\mathcal{S}_q}^-\|_1)$. With $B_3 = \sum_{k=1}^M B_{3k}$ and using (118) and (119), similar to [9, Eqns. (B.60)], we have

$$\begin{aligned} \alpha B_2 + B_3 &\geq -\alpha |B_2| + B_3 \\ &\geq -\alpha \lambda_q \sum_{k=1}^M \|\mathbf{A}_{k\mathcal{S}_q^c}^-\|_1 - \alpha \frac{C_{0q}}{M^{3/2} q} r_{qn} \sum_{k=1}^M \|\mathbf{A}_k\|_F \end{aligned} \quad (122)$$

where we also used the first inequality in (45). Using $\|\mathbf{A}_{k\mathcal{S}_q^c}^-\|_1 \leq \sqrt{s_q} \|\mathbf{A}_k\|_F$, $\sum_{k=1}^M \|\mathbf{A}_k\|_F \leq \sqrt{M} \|\mathbf{A}\|_F$ and the second inequality in (45), we can simplify (122) as

$$\alpha B_2 + B_3 \geq -2\alpha \|\mathbf{A}\|_F \frac{C_{0q}}{Mq} r_{qn}. \quad (123)$$

In a similar manner (see also [9, Eqns. (B.61)]) using (121), we have

$$(1 - \alpha) B_2 + B_4 \geq -2(1 - \alpha) \|\mathbf{A}\|_F \frac{C_{0q}}{Mq} r_{qn} \quad (124)$$

under the upperbound on λ_{qn} specified in (45). Thus, by (112), (117), (123) and (124), we obtain

$$\begin{aligned} J(\mathbf{A}) &\geq \frac{\|\mathbf{A}\|_F^2 \beta_{q,\min}^2}{8Mq} - \frac{2C_{0q} r_{qn} \|\mathbf{A}\|_F}{Mq} \\ &= \frac{\|\mathbf{A}\|_F^2 \beta_{q,\min}^2}{8Mq} \left(1 - \frac{16}{17}\right) > 0 \end{aligned} \quad (125)$$

using $\|\mathbf{A}\|_F = R_q r_{qn}$ and $R_q = 17C_{0q}/\beta_{q,\min}^2$. This proves Theorem 2(ii). \square

APPENDIX C PROOF OF THEOREM 3

Proof of Theorem 3(i). Since $\|\Omega^\circ\|_F = 1$, we have $\bar{\Omega}(\Gamma)/\|\bar{\Omega}(\Gamma)\|_F = \Omega^\circ$. We have

$$\begin{aligned} \|\hat{\Omega} - \Omega^\circ\|_F &= \left\| \frac{\hat{\Omega}(\Gamma)}{\|\hat{\Omega}(\Gamma)\|_F} - \frac{\bar{\Omega}(\Gamma)}{\|\bar{\Omega}(\Gamma)\|_F} \right\|_F \\ &= \left\| \frac{\hat{\Omega}(\Gamma)}{\|\hat{\Omega}(\Gamma)\|_F} - \frac{\bar{\Omega}(\Gamma)}{\|\hat{\Omega}(\Gamma)\|_F} + \frac{\bar{\Omega}(\Gamma)}{\|\hat{\Omega}(\Gamma)\|_F} - \frac{\bar{\Omega}(\Gamma)}{\|\bar{\Omega}(\Gamma)\|_F} \right\|_F \\ &\leq \frac{\|\hat{\Omega}(\Gamma) - \bar{\Omega}(\Gamma)\|_F}{\|\hat{\Omega}(\Gamma)\|_F} + \|\bar{\Omega}(\Gamma)\|_F \left| \frac{1}{\|\hat{\Omega}(\Gamma)\|_F} - \frac{1}{\|\bar{\Omega}(\Gamma)\|_F} \right| \\ &\leq \frac{2}{\|\hat{\Omega}(\Gamma)\|_F} \|\hat{\Omega}(\Gamma) - \bar{\Omega}(\Gamma)\|_F \end{aligned} \quad (126)$$

using $|\|\bar{\Omega}(\Gamma)\|_F - \|\hat{\Omega}(\Gamma)\|_F| \leq \|\hat{\Omega}(\Gamma) - \bar{\Omega}(\Gamma)\|_F$ (by triangle inequality). Now $\|\hat{\Omega}(\Gamma)\|_F = \|\hat{\Omega}(\Gamma) - \bar{\Omega}(\Gamma) + \bar{\Omega}(\Gamma)\|_F \geq \|\bar{\Omega}(\Gamma)\|_F - \|\hat{\Omega}(\Gamma) - \bar{\Omega}(\Gamma)\|_F$. For $n > N_{2p}$, we have $\|\hat{\Omega}(\Gamma) - \bar{\Omega}(\Gamma)\|_F \leq 0.5\|\bar{\Omega}(\Gamma)\|_F$, and therefore, $\|\hat{\Omega}(\Gamma)\|_F \geq 0.5\|\bar{\Omega}(\Gamma)\|_F$. Hence,

$$\|\hat{\Omega} - \Omega^\circ\|_F \leq 4 \|\hat{\Omega}(\Gamma) - \bar{\Omega}(\Gamma)\|_F / \|\bar{\Omega}(\Gamma)\|_F. \quad (127)$$

We now characterize $\|\bar{\Omega}(\Gamma)\|_F$. We have

$$\begin{aligned} A &= \left| \sum_{k=1}^M (\text{tr}(\bar{S}_k^\circ \Phi_k) + \text{tr}(\bar{S}_k^\circ \Phi_k)^*) \right| \leq 2 \sum_{k=1}^M |\text{tr}(\bar{S}_k^\circ \Phi_k)| \\ &\leq 2 \sum_{k=1}^M \|\bar{S}_k^\circ\|_F \|\Phi_k\|_F \leq 2\sqrt{q} \beta_{q,\max} \sum_{k=1}^M \|\Phi_k\|_F. \end{aligned}$$

Since $\sum_{k=1}^M \|\Phi_k\|_F \leq \sum_{k=1}^M \|\Phi_k - \Phi_k^\circ\|_F + \sum_{k=1}^M \|\Phi_k^\circ\|_F \leq \sqrt{M} \gamma_q + \sqrt{q} M / \beta_{q,\min}$, we have $A \leq 0.2\sqrt{q}M + 2qM\beta_{q,\max}/\beta_{q,\min} \leq 2qM(1 + \beta_{q,\max}/\beta_{q,\min})$. By (33) and the fact $\|\Omega^\circ\|_F = 1$, we infer $\|\bar{\Omega}(\Gamma)\|_F \geq \beta_{q,\min}/(\beta_{q,\max} + \beta_{q,\min}) = 1/\gamma_r$, which combined with (127) and (44) yields (47). \square

Proof of Theorem 3(ii). For $n > N_{3p}$, $\hat{\Omega} \in \mathcal{B}(\Omega^\circ)$ (cf. Theorem 3(i)), and $C_{2p} r_{pn} \leq (p/2)$ w.h.p. We have

$$\|\hat{\Gamma}(\hat{\Omega}) - \Gamma^\circ\|_F \leq \|\hat{\Gamma}(\hat{\Omega}) - \bar{\Gamma}(\hat{\Omega})\|_F + \|\bar{\Gamma}(\hat{\Omega}) - \Gamma^\circ\|_F$$

where Theorem 2(ii) applies to $\|\hat{\Gamma}(\hat{\Omega}) - \bar{\Gamma}(\hat{\Omega})\|_F$. By (34),

$$\bar{\Gamma}(\hat{\Omega}) - \Gamma^\circ = \left(\frac{p}{\text{tr}(\Sigma^\circ \hat{\Omega})} - 1 \right) \Gamma^\circ. \quad (128)$$

As in the proof of Theorem 2(ii) (following (116)), we have $\text{tr}(\Sigma^\circ \hat{\Omega}) = \text{tr}(\Sigma^\circ (\hat{\Omega} - \Omega^\circ + \Omega^\circ)) = \text{tr}(\Sigma^\circ (\hat{\Omega} - \Omega^\circ)) + p$ and $|\text{tr}(\Sigma^\circ (\hat{\Omega} - \Omega^\circ))| \leq \|\Sigma^\circ\|_F \|\hat{\Omega} - \Omega^\circ\|_F \leq C_{2p} r_{pn}$ (using (47)). Therefore, $p - C_{2p} r_{pn} \leq \text{tr}(\Sigma^\circ \hat{\Omega}) \leq p + C_{2p} r_{pn}$ and $|p - \text{tr}(\Sigma^\circ \hat{\Omega})| \leq C_{2p} r_{pn}$. Since $0 < C_{2p} r_{pn} \leq (p/2)$ w.h.p., $|\text{tr}(\Sigma^\circ \hat{\Omega})|^{-1} \leq 2/p$. Thus we have $\|\bar{\Gamma}(\hat{\Omega}) - \Gamma^\circ\|_F \leq C_{2q} r_{pn}$, which yields (48). The given probability bound is the result of the bounds in Theorem 2 (both (44) and (46) must hold) and an application of the union bound. \square

REFERENCES

- [1] S.L. Lauritzen, *Graphical Models*. Oxford, UK: Oxford Univ. Press, 1996.
- [2] M.J. Wainwright, *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge, UK: Cambridge Univ. Press, 2019.

- [3] N. Meinshausen and P. Bühlmann, "High-dimensional graphs and variable selection with the Lasso," *Ann. Statist.*, vol. 34, no. 3, pp. 1436-1462, 2006.
- [4] O. Banerjee, L.E. Ghaoui and A. d'Aspremont, "Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data," *J. Mach. Learn. Res.*, vol. 9, pp. 485-516, 2008.
- [5] 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.
- [6] R. Dahlhaus, "Graphical interaction models for multivariate time series," *Metrika*, vol. 51, pp. 157-172, 2000.
- [7] A. Jung, G. Hannak and N. Goertz, "Graphical LASSO based model selection for time series," *IEEE Signal Process. Lett.*, vol. 22, no. 10, pp. 1781-1785, Oct. 2015.
- [8] J.K. Tugnait, "Graphical modeling of high-dimensional time series," in *Proc. 52nd Asilomar Conf. Signals, Systems, Computers*, pp. 840-844, Pacific Grove, CA, Oct. 29 - Oct. 31, 2018.
- [9] J.K. Tugnait, "On sparse high-dimensional graphical model learning for dependent time series," *Signal Process.*, vol. 197, pp. 1-18, Aug. 2022, Article 108539.
- [10] E. Avventi, A. Lindquist, and B. Wahlberg, "ARMA identification of graphical models," *IEEE Trans. Autom. Control*, vol. 58, no. 5, pp. 1167-1178, 2013.
- [11] M. Zorzi and R. Sepulchre, "AR identification of latent-variable graphical models," *IEEE Trans. Autom. Control*, vol. 61, no. 9, pp. 2327-2340, 2016.
- [12] D. Alpagó, M. Zorzi and A. Ferrante, "Identification of sparse reciprocal graphical models," *IEEE Control Sys. Lett.*, vol. 22, no. 4, pp. 659-664, 2018.
- [13] J. Songsiri and L. Vandenberghe, "Topology selection in graphical models of autoregressive processes," *J. Mach. Learn. Res.*, vol. 11, pp. 2671-2705, Oct. 2010.
- [14] V. Ciccone, A. Ferrante and M. Zorzi, "Learning latent variable dynamic graphical models by confidence sets selection," *IEEE Trans. Autom. Control*, vol. 65, no. 12, pp. 5130-5143, Dec. 2020.
- [15] D. Alpagó, M. Zorzi and A. Ferrante, "Data-driven link prediction over graphical models," *IEEE Trans. Autom. Control*, vol. 68, no. 4, pp. 2215-2228, Apr. 2023.
- [16] S. Basu and G. Michailidis, "Regularized estimation in sparse high-dimensional time series models," *Annals Statist.*, vol. 43, no. 4, pp. 1535-1567, 2015.
- [17] C. Leng and C.Y. Tang, "Sparse matrix graphical models," *J. Amer. Statist. Assoc.*, vol. 107, pp. 1187-1200, Sep. 2012.
- [18] 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.
- [19] Y. Zhu and L. Li, "Multiple matrix Gaussian graphs estimation," *J. Royal Statistical Soc., Series B*, vol. 80, pp. 927-950, 2018.
- [20] X. Chen and W. Liu, "Graph estimation for matrix-variate Gaussian data," *Statistica Sinica*, vol. 29, pp. 479-504, 2019.
- [21] K. Greenewald, S. Zhou and A. Hero III, "Tensor graphical lasso (teralasso)," *J. Royal Statistical Soc., Series B*, vol. 81, no. 5, pp. 901-931, 2019.
- [22] 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 Mach. Intell.*, vol. 42, no. 8, pp. 2024-2037, 1 Aug. 2020.
- [23] F. Huang and S. Chen, "Joint learning of multiple sparse matrix Gaussian graphical models," *IEEE Trans. Neural Netw. Learning Sys.*, vol. 26, no. 11, pp. 2606-2620, Nov. 2015.
- [24] S. Zhou, "Gemini: Graph estimation with matrix variate normal instances," *Annals Statist.*, vol. 42, no. 2, pp. 532-562, 2014.
- [25] J. Yin and H. Li, "Model selection and estimation in the matrix normal graphical model," *J. Multivariate Analysis*, vol. 107, pp. 119-140, May 2012.
- [26] S. He, J. Yin, H. Li and X. Wang, "Graphical model selection and estimation for high dimensional tensor data," *J. Multivariate Analysis*, vol. 128, pp. 165-185, 2014.
- [27] K. Min, Q. Mai and X. Zhang, "Fast and separable estimation in high-dimensional tensor Gaussian graphical models," *J. Comp. Graphical Statistics*, vol. 31, pp. 294-300, 2022.
- [28] K. Werner, M. Jansson and P. Stoica, "On estimation of covariance matrices with Kronecker product structure," *IEEE Trans. Signal Process.*, vol. 56, no. 2, pp. 478-491, Feb. 2008.
- [29] P.M. Weichsel, "The Kronecker product of graphs," *Proc. American Math. Soc.*, vol. 13, no. 1, pp. 37-52, 1962.
- [30] J. Leskovec, D. Chakrabarti, J. Kleinberg, C. Faloutsos, and Z. Ghahramani, "Kronecker graphs: An approach to modeling networks," *J. Mach. Learn. Res.*, vol. 11, pp. 985-1042, Feb. 2010.
- [31] C.M. Carvalho and M. West, "Dynamic matrix-variate graphical models," *Bayesian Analysis*, vol. 2, no. 1, pp. 69-98, 2007.
- [32] H. Wang and M. West, "Bayesian analysis of matrix normal graphical models," *Biometrika*, vol. 96, no. 4, pp. 821-834, Dec. 2009.
- [33] Y. Jiang, J. Bigot and S. Maabout, "Online graph topology learning from matrix-valued time series," arXiv:2107.08020v1 [stat.ML], July 2021.
- [34] M. Zorzi, "Nonparametric identification of Kronecker networks," *Automatica*, vol. 145, no. 9, p. 110518, Nov. 2022.
- [35] B. Sinquin and M. Verhaegen, "Quarks: Identification of large-scale Kronecker vector-autoregressive models," *IEEE Trans. Autom. Control*, vol. 64, no. 3, pp. 448-463, 2019.
- [36] M. Zorzi, "Autoregressive identification of Kronecker graphical models," *Automatica*, vol. 119, no. 9, p. 109053, Sep. 2020.
- [37] R. Chen, H. Xiao and D. Yang, "Autoregressive models for matrix-valued time series," *J. Econometrics*, vol. 222, pp. 539-560, 2021.
- [38] P.J. Schreier and L.L. Scharf, *Statistical Signal Processing of Complex-Valued Data*, Cambridge, UK: Cambridge Univ. Press, 2010.
- [39] J.K. Tugnait, "Sparse high-dimensional matrix-valued graphical model learning from dependent data," in *Proc. 22nd IEEE Statistical Signal Process. Workshop (SSP-2023)*, pp. 344-348, Hanoi, Vietnam, July 2-5, 2023.
- [40] J.K. Tugnait, "Sparse-group lasso for graph learning from multi-attribute data," *IEEE Trans. Signal Process.*, vol. 69, pp. 1771-1786, 2021. (Corrections: vol. 69, p. 4758, 2021.)
- [41] A.J. Rothman, P.J. Bickel, E. Levina and J. Zhu, "Sparse permutation invariant covariance estimation," *Elec. J. Statistics*, vol. 2, pp. 494-515, 2008.
- [42] A.K. Gupta and D.K. Nagar, *Matrix Variate Distributions*. Boca Raton, FL: Chapman and Hall/CRC Press, 1999.
- [43] J.K. Tugnait, "Edge exclusion tests for graphical model selection: Complex Gaussian vectors and time series," *IEEE Trans. Signal Process.*, vol. 67, no. 19, pp. 5062-5077, Oct. 1, 2019.
- [44] D.R. Brillinger, *Time Series: Data Analysis and Theory*, Expanded edition. New York: McGraw Hill, 1981.
- [45] 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.
- [46] N. Simon, J. Friedman, T. Hastie and R. Tibshirani, "A sparse-group lasso," *J. Comp. Graphical Statistics*, vol. 22, pp. 231-245, 2013.
- [47] M. Yuan and Y. Lin, "Model selection and estimation in regression with grouped variables," *J. Royal Statistical Soc., Series B*, vol. 68, pp. 49-67, 2006.
- [48] J. Gorski, F. Pfeuffer and K. Klamroth, "Biconvex sets and optimization with biconvex functions: A survey and extensions," *Math. Methods Operations Res.*, vol. 66, pp. 373-408, 2007.
- [49] 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.
- [50] S. Zhang, B. Guo, A. Dong, J. He, Z. Xu and S.X. Chen, "Cautionary tales on air-quality improvement in Beijing," *Proc. Royal Soc. A*, vol. 473, p. 20170457, 2017.
- [51] 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.
- [52] J. Fan, Y. Feng and Y. Wu, "Network exploration via the adaptive lasso and SCAD penalties," *Annals Applied Statistics*, vol. 3, no. 2, pp. 521-541, 2009.
- [53] J.K. Tugnait, "Sparse-group log-sum penalized graphical model learning for time series," in *Proc. 2022 IEEE Intern. Conf. Acoustics, Speech, Signal Process. (ICASSP 2022)*, pp. 5822-5826, Singapore, May 22-27, 2022.