

# SPARSE HIGH-DIMENSIONAL MATRIX-VALUED GRAPHICAL MODEL LEARNING FROM DEPENDENT DATA

Jitendra K. Tugnait

Department of Electrical & Computer Engineering  
Auburn University, Auburn, AL 36849, USA

## 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 matrix graphical models assume that i.i.d. observations of 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 results also yields a rate of convergence. We illustrate our approach using numerical examples.

**Keywords:** Sparse graph learning; matrix graph estimation; matrix time series; undirected graph; inverse spectral density estimation.

## 1. 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  $\mathbf{x}$  has positive-definite covariance matrix  $\mathbf{\Sigma}$  with precision matrix  $\mathbf{\Omega} = \mathbf{\Sigma}^{-1}$ . Then  $\Omega_{ij}$ , the  $(i, j)$ -th element of  $\mathbf{\Omega}$ , is zero iff  $x_i$  and  $x_j$  are conditionally independent [1]. Such models for  $\mathbf{x}$  have been extensively studied [2–4].

These models are vector graphical models. Time series (dependent data) graphical models are much less studied. 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  $\{x_\ell(t), \ell \in [p], \ell \neq i, \ell \neq j\}$  [5]. Vector graphical models (based only on the precision matrix  $\mathbf{\Omega}$ ) do not necessarily capture the “true” series graphical model if the data originates from a time-dependent series. Denote the power spectral density (PSD) matrix of  $\{\mathbf{x}(t)\}$  by  $\mathbf{S}_x(f)$ , where  $\mathbf{S}_x(f) = \sum_{\tau=-\infty}^{\infty} \mathbf{R}_{xx}(\tau) e^{-j2\pi f\tau}$ ,  $\mathbf{R}_{xx}(\tau) = E\{\mathbf{x}(t+\tau)\mathbf{x}^\top(t)\}$ . In [5] it was shown that conditional independence of two time series components given all other components of the time series, is encoded by

zeros in the inverse PSD, that is,  $\{i, j\} \notin \mathcal{E}$  iff the  $(i, j)$ -th element of  $\mathbf{S}_x(f)$ ,  $[\mathbf{S}_x^{-1}(f)]_{ij} = 0$  for every  $f$ .

The need for matrix-valued graphical models arises in several applications [6–15]. Here we observe matrix-valued time series  $\{\mathbf{Z}(t)\}$  where  $\mathbf{Z}(t) \in \mathbb{R}^{p \times q}$ . If one vectorizes using  $\text{vec}(\mathbf{Z})$ , 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 moreover, it ignores any structural information among rows and columns of the matrix observations [6]. Prior work [6–15] all assume that i.i.d. observations of  $\mathbf{Z}$  are available for graphical modeling. Our objective in this paper is to learn the graph associated with time-dependent matrix-valued  $p \times q$  Gaussian sequence  $\mathbf{Z}(t)$ , given observations of  $\mathbf{Z}(t)$  for  $t = 0, 1, \dots, n-1$ .

*Notation:*  $|\mathbf{A}|$  and  $\text{tr}(\mathbf{A})$  denote the determinant and the trace of the square matrix  $\mathbf{A}$ , respectively, and  $\text{etr}(\mathbf{A}) = \exp(\text{tr}(\mathbf{A}))$ .  $[\mathbf{B}]_{ij}$  denotes the  $(i, j)$ -th element of  $\mathbf{B}$ , and so does  $B_{ij}$ .  $\mathbf{I}_m$  is the  $m \times m$  identity matrix. The superscripts  $*$  and  $H$  denote the complex conjugate and conjugate transpose operations, respectively,  $\mathbf{x} \sim \mathcal{N}_c(\mathbf{m}, \mathbf{\Sigma})$  denotes a random vector  $\mathbf{x}$  that is circularly symmetric (proper) complex Gaussian with mean  $\mathbf{m}$  and covariance  $\mathbf{\Sigma}$ , and  $\otimes$  denotes the Kronecker product.

## 2. SYSTEM MODEL

Random matrix  $\mathbf{Z} \in \mathbb{R}^{p \times q}$  is said to have a matrix normal (Gaussian) distribution if its pdf  $f(\mathbf{Z}|\mathbf{M}, \mathbf{\Sigma}, \mathbf{\Psi})$ , characterized by  $\mathbf{M} \in \mathbb{R}^{p \times q}$ ,  $\mathbf{\Sigma} \in \mathbb{R}^{p \times p}$ ,  $\mathbf{\Psi} \in \mathbb{R}^{q \times q}$ , is given by [16, Chap. 2]

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

Equivalently,

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

Here  $\mathbf{\Psi}$  is the row covariance matrix and  $\mathbf{\Sigma}$  is the column covariance matrix [16] since the  $k$ th column  $\mathbf{Z}_{\cdot k} \sim \mathcal{N}(\mathbf{0}, [\mathbf{\Psi}]_{kk} \mathbf{\Sigma})$  and the  $i$ th row  $\mathbf{Z}_i^\top \sim \mathcal{N}(\mathbf{0}, [\mathbf{\Sigma}]_{ii} \mathbf{\Psi})$ .

Graphical modeling of random vectors to characterize conditional dependence of its components [1, 3] has been extended to matrix data with structured information [6–9, 13]. With  $\mathbf{Z} \in \mathbb{R}^{p \times q}$  modeled as a zero-mean matrix normal vector and  $\mathbf{z} = \text{vec}(\mathbf{Z})$ , [6] assumes

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

which could be interpreted as follows. Let  $\mathbf{x} \in \mathbb{R}^p$ ,  $\mathbf{y} \in \mathbb{R}^q$ ,  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})$  and  $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Psi})$ . Express  $\mathbf{Z}$  as

$$\mathbf{Z} = \mathbf{x} \otimes \mathbf{y}^\top, \text{ or } \mathbf{z} = \text{vec}(\mathbf{Z}) = \text{vec}(\mathbf{x}\mathbf{y}^\top) = \mathbf{y} \otimes \mathbf{x} \quad (4)$$

This work was supported by NSF Grant ECCS-2040536. Author's email: tugnajk@auburn.edu

such that

$$\begin{aligned} E\{\mathbf{z}\mathbf{z}^\top\} &= E\{(\mathbf{y} \otimes \mathbf{x})(\mathbf{y} \otimes \mathbf{x})^\top\} = E\{(\mathbf{y}\mathbf{y}^\top) \otimes (\mathbf{x}\mathbf{x})^\top\} \\ &= E\{\mathbf{y}\mathbf{y}^\top\} \otimes E\{\mathbf{x}\mathbf{x}^\top\} = \Psi \otimes \Sigma, \end{aligned} \quad (5)$$

implying a separable covariance structure [17]. Let  $\Omega = \Sigma^{-1}$  and  $\Gamma = \Psi^{-1}$  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  $\Omega_{ij}$  and  $\Gamma_{k\ell}$  is zero when  $i \neq k, j \neq \ell$ , (ii)  $\Omega_{ij} = 0$  when  $i \neq k, j = \ell$ , and (iii)  $\Gamma_{k\ell} = 0$  when  $i = k, j \neq \ell$  [6]. Prior work [6–9, 13] all assume that i.i.d. observations of  $\mathbf{Z}$  are available for graphical modeling.

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

$$E\{\mathbf{z}(t + \tau)\mathbf{z}^\top(t)\} = \Psi(\tau) \otimes \Sigma \quad (6)$$

where  $\Psi(\tau)$ ,  $\tau = 0, \pm 1, \dots$  models time-dependence while  $\Sigma \succ \mathbf{0}$  is fixed. With  $\{\mathbf{e}(t)\}$  i.i.d.,  $\mathbf{e}(t) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , a generative model for  $\mathbf{z}(t)$  is given by

$$\begin{aligned} \mathbf{z}(t) &= \sum_{i=0}^L (\mathbf{B}_i \otimes \mathbf{F}) \mathbf{e}(t - i), \quad \mathbf{B}_i \in \mathbb{R}^q, \quad \mathbf{F} \in \mathbb{R}^p \quad (7) \\ \Rightarrow E\{\mathbf{z}(t + \tau)\mathbf{z}^\top(t)\} &= \underbrace{\left( \sum_{i=0}^L \mathbf{B}_i \mathbf{B}_{i-\tau}^\top \right)}_{=\Psi(\tau)} \otimes \underbrace{(\mathbf{F}\mathbf{F}^\top)}_{=\Sigma}. \quad (8) \end{aligned}$$

An example considered in [6] is that of a United States Department of Agriculture (USDA) dataset reporting itemized annual export to major trading partners. The dataset with 40 years U.S. export is collected for 13 trading partners and 36 items. Each observation in the dataset can be denoted by a  $13 \times 36$  matrix where the trading partners and items, as the rows and columns, respectively, of this matrix, are used as structural information for the observations. The basic idea in matrix-valued graphs is to model the covariance of  $\text{vec}(\mathbf{Z})$  as  $\Psi \otimes \Sigma$  reducing the number of unknowns from  $\mathcal{O}(p^2q^2)$  in the precision matrix for the “full” vectorized model to  $\mathcal{O}(p^2 + q^2)$  for the matrix model, while also preserving the structural information. Given data, one estimates two precision matrices  $\Omega = \Sigma^{-1}$  and  $\Gamma = \Psi^{-1}$ . In the matrix graph, conditional independence between  $\mathbf{Z}_{ij}$  and  $\mathbf{Z}_{k\ell}$  is determined by zeros in  $\Omega$  and  $\Gamma$  [6]. While [6] and others ([8–11]) all consider only i.i.d. observations, we allow possible temporal dependence in matrix observations via  $\Psi(\tau)$ .

The PSD of  $\{\mathbf{z}(t)\}$  is  $\mathbf{S}_z(f) = \bar{\mathbf{S}}(f) \otimes \Sigma$  where  $\bar{\mathbf{S}}(f) = \sum_{\tau} \Psi(\tau) e^{-j2\pi f\tau}$ . Then  $\mathbf{S}_z^{-1}(f) = \bar{\mathbf{S}}^{-1}(f) \otimes \Sigma^{-1}$ , and by [5], in the  $pq$ -node graph  $\mathcal{G} = (V, \mathcal{E})$ ,  $|V| = pq$ , associated with  $\{\mathbf{z}(t)\}$ , edge  $\{i, j\} \in \mathcal{E}$  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  $\Gamma = \Psi^{-1}$ , using [5, 6], 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  $\Omega_{ij}$  and  $[\bar{\mathbf{S}}^{-1}(f)]_{k\ell}$ ,  $f \in [0, 0.5]$  is zero when  $i \neq k, j \neq \ell$ , (ii)  $\Omega_{ij} = 0$  when  $i \neq k, j = \ell$ , and (iii)  $[\bar{\mathbf{S}}^{-1}(f)]_{k\ell} = 0$  for  $f \in [0, 0.5]$  when  $i = k, j \neq \ell$ .

As an example, consider  $\mathbf{x}(t) = \sum_{i=0}^L \mathbf{B}_i \mathbf{e}(t - i)$ ,  $L \geq 1$ , where  $\{\mathbf{e}(t)\}$  is zero-mean, i.i.d. Gaussian, with covariance  $= \mathbf{I}$ ,

$$\begin{aligned} \mathbf{B}_0 &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & \sqrt{1-b^2} & 0 \\ a & 0 & \sqrt{1-a^2} \end{bmatrix}, \quad \mathbf{B}_L = \begin{bmatrix} 0 & 0 & 0 \\ b & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \\ \mathbf{B}_i &= \mathbf{0} \quad \forall i \neq 0, i \neq L, \quad |a| < 1, |b| < 1. \end{aligned} \quad (9)$$

Straightforward calculations yield

$$\Gamma_x = \mathbf{R}_{xx}^{-1}(0) = \frac{1}{1-a^2} \begin{bmatrix} 1 & 0 & -a \\ 0 & 1-a^2 & 0 \\ -a & 0 & 1 \end{bmatrix}, \quad (10)$$

$$c \mathbf{S}_x^{-1}(f) = \begin{bmatrix} 1 - (ab)^2 & d_1 & -a(1-b^2) \\ d_1^* & 1-a^2 & 0 \\ -a(1-b^2) & 0 & 1-b^2 \end{bmatrix} \quad (11)$$

where  $c = (1-a^2)(1-b^2)$  and  $d_1 = -b(1-a^2)e^{j2\pi fL}$ . Notice that in  $\Gamma_x$ , edges  $\{1, 2\}$  and  $\{2, 1\}$  are missing whereas they are present in  $\mathbf{S}_x^{-1}(f)$ , that is,  $\Gamma_x$  does not capture the true dependencies among various components of the dependent series.

Our objective is to learn the graph associated with time-dependent sequence  $\{\mathbf{Z}(t)\}$ , given observations  $t = 0, 1, \dots, n-1$ , under some sparsity constraints on  $\Omega$  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 will take  $\Omega_{11} = 1$ .

### 3. PENALIZED NEGATIVE LOG-LIKELIHOOD

Given  $\mathbf{z}(t)$  for  $t = 0, 1, 2, \dots, n-1$ . Define the (normalized) DFT  $\mathbf{d}_z(f_m)$  of  $\mathbf{z}(t)$ , ( $j = \sqrt{-1}$ ,  $f_m = m/n$ ), over  $m = 0, 1, \dots, n-1$  as  $\mathbf{d}_z(f_m) = \frac{1}{\sqrt{n}} \sum_{t=0}^{n-1} \mathbf{z}(t) \exp(-j2\pi f_m t)$ . Let  $\mathbf{D}_z(f_m) = \frac{1}{\sqrt{n}} \sum_{t=0}^{n-1} \mathbf{Z}(t) \exp(-j2\pi f_m t)$ , then  $\mathbf{d}_z(f_m) = \text{vec}(\mathbf{D}_z(f_m))$ . It is established in [18] (see also [19]) that 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; in our case, this assumption applies to  $\bar{\mathbf{S}}(f_k)$ . Pick  $M = \lfloor (\frac{n}{2} - m_t - 1)/K \rfloor$  and

$$\tilde{f}_k = \frac{(k-1)K + m_t + 1}{n}, \quad k = 1, 2, \dots, M, \quad (12)$$

yielding  $M$  equally spaced frequencies  $\tilde{f}_k$  in the interval  $(0, 0.5)$ . By local smoothness

$$\mathbf{S}_z(\tilde{f}_{k,\ell}) = \mathbf{S}_z(\tilde{f}_k) \quad \text{for } \ell = -m_t, -m_t + 1, \dots, m_t, \quad (13)$$

$$\text{where } \tilde{f}_{k,\ell} = \frac{(k-1)K + m_t + 1 + \ell}{n}. \quad (14)$$

It is known ([20, Theorem 4.4.1]) that 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, provided all elements of  $\mathbf{R}_{zz}(\tau) = E\{\mathbf{z}(t + \tau)\mathbf{z}^\top(t)\}$  are absolutely summable. Denote the joint probability density function of  $\mathbf{d}_z(f_m)$ ,  $m = 1, 2, \dots, (n/2) - 1$ , as  $f_{\mathcal{D}}$ . Then we have [18, 19]

$$f_{\mathcal{D}}(\mathcal{D}) = \prod_{k=1}^M \left[ \prod_{\ell=-m_t}^{m_t} \frac{\exp(-g - g^*)}{\pi^{pq} |\bar{\mathbf{S}}(\tilde{f}_k) \otimes \Sigma|^{1/2} |\bar{\mathbf{S}}^*(\tilde{f}_k) \otimes \Sigma|^{1/2}} \right] \quad (15)$$

$$\text{where } g = \frac{1}{2} \mathbf{d}_z^H(\tilde{f}_{k,\ell}) (\bar{\mathbf{S}}^{-1}(\tilde{f}_k) \otimes \Sigma^{-1}) \mathbf{d}_z(\tilde{f}_{k,\ell}). \quad (16)$$

Using  $\text{tr}(\mathbf{A}^\top \mathbf{B} \mathbf{C} \mathbf{D}^\top) = (\text{vec}(\mathbf{A}))^\top (\mathbf{D} \otimes \mathbf{B}) \text{vec}(\mathbf{C})$  and  $|\bar{\mathbf{S}}(\tilde{f}_k) \otimes \Sigma| = |\bar{\mathbf{S}}(\tilde{f}_k)|^p |\Sigma|^q$ , and parametrizing in terms of  $\Phi_k := \bar{\mathbf{S}}^{-1}(\tilde{f}_k)$  and  $\Omega = \Sigma^{-1}$ , up to some constants the negative log-likelihood follows from (15) as ( $\{\Phi\}$  denotes  $\{\Phi_k, k =$

$1, \dots, M\}$ )

$$-\frac{1}{KMpq} \ln f_{\mathcal{D}}(\mathcal{D}) \propto G(\Omega, \{\Phi\}, \{\Phi^*\}) := -\frac{1}{p} \ln(|\Omega|) - \frac{1}{2Mq} \sum_{k=1}^M (\ln |\Phi_k| + \ln |\Phi_k^*| - \text{tr}(\mathbf{A}_k + \mathbf{A}_k^*)) \quad (17)$$

$$\text{where } \mathbf{A}_k = \frac{1}{Kp} \sum_{\ell=-m_t}^{m_t} \mathbf{D}_z^H(\tilde{f}_{k,\ell}) \Omega \mathbf{D}_z(\tilde{f}_{k,\ell}) \Phi_k^*. \quad (18)$$

In the high-dimension case, one needs to use penalty terms to enforce sparsity and to make the problem well-conditioned. Imposing a sparse-group lasso sparsity constraint on  $\{\Phi\}$  (cf. [2, 21, 22]) and a lasso constraint on  $\Omega$ , we propose to minimize a penalized version of negative log-likelihood w.r.t.  $\Omega$  and  $\{\Phi\}$ ,

$$\mathcal{L}(\Omega, \{\Phi\}) = G(\Omega, \{\Phi\}, \{\Phi^*\}) + P_p(\Omega) + P_q(\{\Phi\}), \quad (19)$$

$$P_p(\Omega) = \lambda_p \sum_{i \neq j}^p |\Omega_{ij}| \quad (20)$$

$$P_q(\{\Phi\}) = \alpha \lambda_q \sum_{k=1}^M \sum_{i \neq j}^p |[\Phi_k]_{ij}| + (1 - \alpha) \sqrt{M} \lambda_q \sum_{i \neq j}^p \|\Phi^{(ij)}\| \quad (21)$$

$$\text{where } \Phi^{(ij)} := [[\Phi_1]_{ij} [\Phi_2]_{ij} \dots [\Phi_M]_{ij}]^\top \in \mathbb{C}^M, \quad (22)$$

and  $\alpha \in [0, 1]$  and  $\lambda_p, \lambda_q > 0$  are tuning parameters.

#### 4. OPTIMIZATION

The objective function  $\mathcal{L}(\Omega, \{\Phi\})$  in (19) is biconvex: (strictly) convex in  $\{\Phi\}$ ,  $\Phi_k \succ \mathbf{0}$ , for fixed  $\Omega$ , and (strictly) convex in  $\Omega$ ,  $\Omega \succ \mathbf{0}$ , for fixed  $\{\Phi\}$ . As in [6, 7] (and others) pertaining to the i.i.d. observations case, and as is a general approach for biconvex function optimization [23], we will use an iterative and alternating minimization approach where we optimize w.r.t.  $\Omega$  with  $\{\Phi\}$  fixed, and then optimize w.r.t.  $\{\Phi\}$  with  $\Omega$  fixed at the last optimized value, and repeat the two optimizations (flip-flop). There is no guarantee that the algorithm converges to the global minimum, however, the algorithm converges to a local stationary point of  $\mathcal{L}(\Omega, \{\Phi\})$  [23].

With  $\{\hat{\Phi}\}$  denoting the estimate of  $\{\Phi\}$ , fix  $\{\Phi\} = \{\hat{\Phi}\}$  and let  $\mathcal{L}_1(\Omega)$  denote  $\mathcal{L}(\Omega, \{\hat{\Phi}\})$  up to some irrelevant constants. We minimize  $\mathcal{L}_1(\Omega)$  w.r.t.  $\Omega$  to obtain estimate  $\hat{\Omega}$ , where

$$\mathcal{L}_1(\Omega) = -\frac{1}{p} \ln(|\Omega|) + \frac{1}{p} \text{tr}(\Omega \bar{\mathbf{S}}) + P_p(\Omega), \quad (23)$$

$$\bar{\mathbf{S}} = \frac{1}{MKq} \sum_{k=1}^M \sum_{\ell=-m_t}^{m_t} \text{Re}\{\mathbf{D}_z(\tilde{f}_{k,\ell}) \hat{\Phi}_k^* \mathbf{D}_z^H(\tilde{f}_{k,\ell})\}. \quad (24)$$

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

$$\mathcal{L}_2(\{\Phi\}) = -\frac{1}{2Mq} \sum_{k=1}^M (\ln |\Phi_k| + \ln |\Phi_k^*|) + \frac{1}{2Mq} \sum_{k=1}^M \text{tr}(\tilde{\mathbf{S}}_k \Phi_k + \tilde{\mathbf{S}}_k^* \Phi_k^*) + P_q(\{\Phi\}), \quad (25)$$

$$\tilde{\mathbf{S}}_k = \frac{1}{Kp} \sum_{\ell=-m_t}^{m_t} \mathbf{D}_z^\top(\tilde{f}_{k,\ell}) \hat{\Omega} \mathbf{D}_z^*(\tilde{f}_{k,\ell}). \quad (26)$$

Our optimization algorithm (used in our simulations) is as follows.

1. Initialize  $m = 1$ ,  $\Omega^{(0)} = \mathbf{I}_p$ ,  $\Phi_k^{(0)} = \mathbf{I}_q$ ,  $k = 1, 2, \dots, M$ .
2. Set  $\hat{\Omega} = \Omega^{(m-1)}$  in (26). Use the iterative alternating direction method of multipliers (ADMM) algorithm [24], as outlined in [19, Sec. 4], to minimize  $\mathcal{L}_2(\{\Phi\})$  (given by (25)) w.r.t.  $\{\Phi\}$  to obtain estimates  $\Phi_k^{(m)}$ ,  $k = 1, 2, \dots, M$ . Cost (40) in [19] corresponds to (25) of this paper.
3. Set  $\{\hat{\Phi}\} = \{\hat{\Phi}^{(m)}\}$  in (24). Use the ADMM algorithm of [25, Sec. III] (with  $\alpha = 1$  therein, no group-lasso penalty) to minimize  $\mathcal{L}_1(\Omega)$  w.r.t.  $\Omega$ , to obtain estimate  $\Omega^{(m)}$ . Cost (7) in [25] (after setting  $\alpha = 1$ ) corresponds to (23) of this paper. Normalize  $\hat{\Omega}_{11}^{(m)} = 1$  to resolve the scaling ambiguity. Let  $m \leftarrow m + 1$ .
4. Repeat steps 2 and 3 until convergence.

#### 4.1. BIC for selection of $\lambda_p, \lambda_q$ (and $\alpha$ )

Given  $n, K$  and  $M$ , the Bayesian information criterion (BIC) is given by (see also [19])  $\text{BIC}(\lambda_p, \lambda_q, \alpha) = -2KMq \ln |\hat{\Omega}| + 2Kp \sum_{k=1}^M (-\ln |\hat{\Phi}_k| + p^{-1} \text{tr}(\hat{\mathbf{A}}_k)) + \ln(2KM) (|\hat{\Omega}|_0/2 + \sum_{k=1}^M |\hat{\Phi}_k|_0)$  where  $\hat{\mathbf{A}}_k$  is given by (18) with  $\Omega$  and  $\Phi_k$  therein replaced with  $\hat{\Omega}$  and  $\hat{\Phi}_k$ , respectively,  $|\mathbf{J}|_0$  denotes number of nonzero elements in  $\mathbf{J}$ ,  $2KM$  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  $(0, \pi)$ . Pick  $\alpha, \lambda_q$  and  $\lambda_p$  to minimize BIC. In our simulations we fixed  $\alpha = 0.05$  and then picked  $\lambda_q$  and  $\lambda_p$  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 [25]. Find the smallest  $\lambda_q$  and  $\lambda_p$ , labeled  $\lambda_{qsm}$  and  $\lambda_{psm}$ , 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  $\lambda_{pu}$  and  $\lambda_{p\ell}$ .

#### 5. CONSISTENCY

Now we provide sufficient conditions for local convergence in the Frobenius norm of the Kronecker-decomposable inverse PSD estimators to the true value. Define  $q \times (qM)$  matrix  $\tilde{\Omega}$  as

$$\tilde{\Omega} = [\Phi_1 \Phi_2 \dots \Phi_M]. \quad (27)$$

We now allow  $p, q, M, K$  (see (12)),  $\lambda_p$  and  $\lambda_q$  to be functions of sample size  $n$ , denoted as  $p_n, q_n, M_n, K_n, \lambda_{pn}$  and  $\lambda_{qn}$ , respectively. Assume

- (A1) 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  $k, \ell \in [q]$ .
- (A2) Define the true edgesets  $\mathcal{S}_q = \{\{i, j\} : [\bar{\mathbf{S}}_0^{-1}(f)]_{ij} \neq 0, i \neq j, 0 \leq f \leq 0.5, i, j \in [q]\}$  and  $\mathcal{S}_p = \{\{i, j\} : \Omega_{ij} \neq 0, i \neq j, i, j \in [p]\}$ , where  $\bar{\mathbf{S}}_0(f)$  denotes DTFT of  $\Psi(\tau)$  and  $\Omega_0 = \Sigma_0^{-1}$  denotes the true value of  $\Omega$ . Assume that  $|\mathcal{S}_q| \leq s_{qn}$  and  $|\mathcal{S}_p| \leq s_{pn}$ .
- (A3) The minimum and maximum eigenvalues of  $q_n \times q_n$  PSD  $\bar{\mathbf{S}}_0(f) \succ \mathbf{0}$  satisfy  $0 < \beta_{q,\min} \leq \min_{f \in [0,0.5]} \phi_{\min}(\bar{\mathbf{S}}_0(f))$  and  $\max_{f \in [0,0.5]} \phi_{\max}(\bar{\mathbf{S}}_0(f)) \leq \beta_{q,\max} < \infty$ . Similarly,  $0 < \beta_{p,\min} \leq \phi_{\min}(\Sigma_0) \leq \phi_{\max}(\Sigma_0) \leq \beta_{p,\max} < \infty$ . Here  $\beta_{\min}$  and  $\beta_{\max}$  are not functions of  $n$  (or  $p_n, q_n$ ).

Theorem 1 whose proof is omitted for lack of space, establishes consistency of a local minimizer  $(\hat{\Omega}, \hat{\Phi})$  of  $\mathcal{L}(\Omega, \{\Phi\}) = \mathcal{L}(\Omega, \tilde{\Omega})$

under assumptions (A1)-(A3). First we define some variables. For  $\tau > 2$ , define

$$C_{0q} = 80 \max_{\ell, f}([\bar{\mathbf{S}}_0(f)]_{\ell\ell}) \sqrt{2 \ln(16q_n^\tau M_n) / \ln(q_n)}, \quad (28)$$

$$C_{0p} = 40 \max_k([\mathbf{\Sigma}_0]_{\ell\ell}) \sqrt{2 \ln(4p_n^\tau) / \ln(p_n)}, \quad (29)$$

$$r_{qn} = \sqrt{M_n(q_n + s_{qn}) \ln(q_n) / (K_n p_n)} = o(1), \quad (30)$$

$$r_{pn} = \sqrt{(p_n + s_{pn}) \ln(p_n) / (M_n K_n q_n)} = o(1). \quad (31)$$

Recall that for random vectors  $\mathbf{y}_n$  and  $\mathbf{x}_n$ , the notation  $\mathbf{y}_n = \mathcal{O}_P(\mathbf{x}_n)$  means that for any  $\varepsilon > 0$ , there exist real  $R$  and integer  $N$ ,  $0 < R < \infty$  and  $0 < N < \infty$ , such that  $P(\|\mathbf{y}_n\| \leq R\|\mathbf{x}_n\|) \geq 1 - \varepsilon \forall n \geq N$ .

*Theorem 1 (Consistency).* For  $\alpha \in [0, 1]$ , any  $C_1 \geq 1$ , and  $C_{0p}$  and  $C_{0q}$  as defined in (28) and (29), respectively, suppose  $\lambda_{pn}$  and  $\lambda_{qn}$  satisfy

$$\frac{C_{0p}}{p_n} \sqrt{\frac{\ln(p_n)}{M_n K_n q_n}} \leq \lambda_{pn} \leq \frac{C_1 C_{0p}}{p_n} \sqrt{\left(1 + \frac{p_n}{s_{pn}}\right) \frac{\ln(p_n)}{M_n K_n q_n}},$$

$$\frac{C_{0q}}{M_n q_n} \sqrt{\frac{\ln(q_n)}{K_n p_n}} \leq \lambda_{qn} \leq \frac{C_1 C_{0q}}{M_n q_n} \sqrt{\left(1 + \frac{q_n}{s_{qn}}\right) \frac{\ln(q_n)}{K_n p_n}}.$$

Then under assumptions (A1)-(A3), there exists a local minimizer  $(\hat{\mathbf{\Omega}}, \hat{\mathbf{\Gamma}})$  of  $\mathcal{L}(\mathbf{\Omega}, \mathbf{\Gamma})$  such that

$$\|\hat{\mathbf{\Omega}} - \mathbf{\Omega}_0\|_F = \mathcal{O}_P(r_{pn}), \quad \|\hat{\mathbf{\Gamma}} - \mathbf{\Gamma}_0\|_F = \mathcal{O}_P(r_{qn}) \quad (32)$$

where  $r_{pn}$  and  $r_{qn}$  are as in (30) and (31), respectively. •

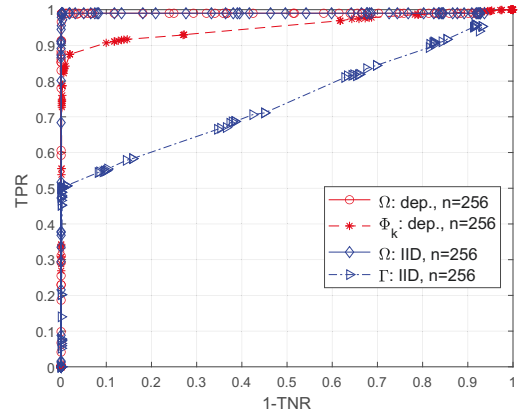
**Remark 1.** Proof of Theorem 1 is patterned after [11] pertaining to matrix graphs, exploiting the results in [19] for dependent vector time series and in [25] for multi-attribute graphical models with i.i.d. data; in turn, all these results are based on the proof technique of [26]. Theorem 1 helps determine how to choose  $M_n$  and  $K_n$  so that for given  $n$ ,  $q_n$  and  $p_n$ ,  $\lim_{n \rightarrow \infty} r_{pn} = 0$  and  $\lim_{n \rightarrow \infty} r_{qn} = 0$ . See also [19, Remark 2]. □

## 6. NUMERICAL RESULTS

We use model (7)-(8) to generate synthetic data where  $\Psi(\tau)$  is controlled via  $\mathbf{B}_i$ 's as in (9) and  $\mathbf{\Sigma}$  is determined via an Erdős-Rényi graph. We take  $p = q = 15$  and for  $\Psi(\tau)$ ,  $\mathbf{B}_i$ 's in (7) have a block-diagonal structure with 5 blocks, each block as in (9), where in each  $3 \times 3$  block,  $a, b$  are uniform over  $[-0.3, 0.7]$ ,  $L$  is uniform over  $\{1, 2, 3, 4, 5\}$ . In the Erdős-Rényi graph with  $p = 15$  nodes, the nodes are connected with probability  $p_{er} = 0.05$ . In the upper triangular  $\mathbf{\Omega}$ ,  $\mathbf{\Omega}_{ij} = 0$  if  $\{i, j\} \notin \mathcal{S}_p$ ,  $\mathbf{\Omega}_{ij}$  is uniformly distributed over  $[-0.4, -0.1] \cup [0.1, 0.4]$  if  $\{i, j\} \in \mathcal{S}_p$ , and  $\mathbf{\Omega}_{ii} = 0.5$ . With  $\mathbf{\Omega} = \mathbf{\Omega}^\top$ , add  $\kappa \mathbf{I}$  to  $\mathbf{\Omega}$  with  $\kappa$  picked to make minimum eigenvalue of  $\mathbf{\Omega} = \mathbf{\Omega} + \kappa \mathbf{I}$  equal to 0.5. Let  $\mathbf{\Omega} = \mathbf{\tilde{F}} \mathbf{\tilde{F}}^\top$  (matrix square-root), then  $\mathbf{F} = \mathbf{\tilde{F}}^{-1}$  in (7).

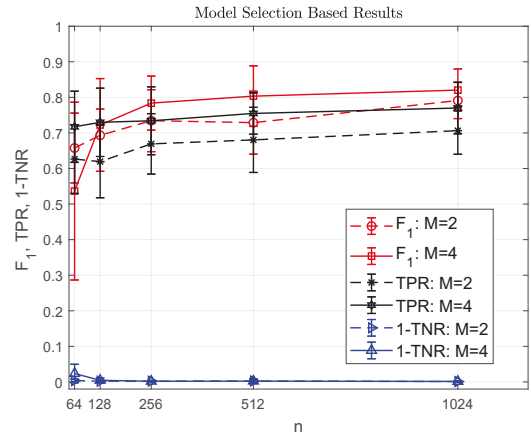
We applied our proposed approach with  $n = 256$ ,  $M = 2$ ,  $K = 63$  and compared with the approach of [6] (which is also the approach of [7, 14], all of whom assume i.i.d. observations and have two lasso penalties one each on  $\mathbf{\Omega}$  and  $\mathbf{\Gamma}$ , counterpart to our  $\mathbf{\Phi}_k$ ). By changing the penalty parameters and determining the resulting edges, we calculated the true positive rate (TPR) and false positive rate 1-TNR (where TNR is the true negative rate) over 100 runs, separately for  $\mathbf{\Omega}$  and  $\{\mathbf{\Phi}_k\}/\mathbf{\Gamma}$ . The receiver operating characteristic (ROC) is shown in Fig. 1 based on 100 runs. Fig. 1 shows that the i.i.d. modeling of [6, 7, 14] is unable to capture the “dependent” edges

(cf. (7)) via  $\mathbf{\Gamma}$  whereas it has no issues with  $\mathbf{\Omega}$ . Our approach works well for both components of the graph Kronecker product.



**Fig. 1:** ROC curves: plots labeled “IID” are from the approach of [6, 7, 14], and the plots labeled “dep.” are from our proposed approach. TPR=true positive rate, TNR=true negative rate

In Fig. 2 we show the results based on 50 runs for our approach when BIC parameter selection method (Sec. 4.1) is applied. We take  $n = 64, 128, 256, 512, 1024$  with corresponding  $m_t$  values as either  $m_t = 7, 15, 31, 63, 127$  ( $M = 2$ ), or  $m_t = 3, 7, 14, 31, 63$  ( $M = 4$ ); note  $K = 2m_t + 1$ . Here we show the TPR, 1-TNR and  $F_1$  score values for the overall graph (not the two Kronecker product components separately) along with the  $\pm \sigma$  error bars. The proposed approach works well both in terms of  $F_1$  score and TPR vs 1-TNR.



**Fig. 2:** BIC based results of the proposed approach:  $F_1$ -scores, TPR and 1-TNR

## 7. CONCLUSIONS

Sparse-group lasso penalized log-likelihood approach in frequency-domain with a Kronecker-decomposable PSD was presented for matrix graph 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 a numerical example where our approach significantly outperformed an existing i.i.d. modeling-based approach [6, 7, 14] in correctly detecting the graph edges with ROC as the performance metric.

## 8. REFERENCES

- [1] S.L. Lauritzen, *Graphical Models*. Oxford, UK: Oxford Univ. Press, 1996.
- [2] P. Danaher, P. Wang and D.M. Witten, "The joint graphical lasso for inverse covariance estimation across multiple classes," *J. Royal Statistical Society, Series B (Methodological)*, vol. 76, pp. 373-397, 2014.
- [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] K. Mohan, P. London, M. Fazel, D. Witten and S.I. Lee, "Node-based learning of multiple Gaussian graphical models," *J. Machine Learning Research*, vol. 15, pp. 445-488, 2014.
- [5] R. Dahlhaus, "Graphical interaction models for multivariate time series," *Metrika*, vol. 51, pp. 157-172, 2000.
- [6] C. Leng and C.Y. Tang, "Sparse matrix graphical models," *J. American Statistical Association*, vol. 107, pp. 1187-1200, Sep. 2012.
- [7] 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.
- [8] Y. Zhu and L. Li, "Multiple matrix Gaussian graphs estimation," *J. Royal Statistical Society, Series B (Methodological)*, vol. 80, pp. 927-950, 2018.
- [9] X. Chen and W. Liu, "Graph estimation for matrix-variate Gaussian data," *Statistica Sinica*, vol. 29, pp. 479-504, 2019.
- [10] K. Greenewald, S. Zhou and A. Hero III, "Tensor graphical lasso (teralasso)," *J. Royal Statistical Society, Series B (Methodological)*, vol. 81, no. 5, pp. 901-931, 2019.
- [11] X. Lyu, W.W. Sun, Z. Wang, H. Liu, J. Yang and G. Cheng, "Tensor graphical model: Non-convex optimization and statistical inference," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 42, no. 8, pp. 2024-2037, 1 Aug. 2020.
- [12] F. Huang and S. Chen, "Joint learning of multiple sparse matrix Gaussian graphical models," *IEEE Trans. Neural Networks and Learning Systems*, vol. 26, no. 11, pp. 2606-2620, Nov. 2015.
- [13] S. Zhou, "Gemini: Graph estimation with matrix variate normal instances," *Annals Statistics*, vol. 42, no. 2, pp. 532-562, 2014.
- [14] 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.
- [15] 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.
- [16] A.K. Gupta and D.K. Nagar, *Matrix Variate Distributions*. Boca Raton, FL: Chapman and Hall/CRC Press, 1999.
- [17] 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.
- [18] 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.
- [19] J.K. Tugnait, "On sparse high-dimensional graphical model learning for dependent time series," *Signal Processing*, vol. 197, pp. 1-18, Aug. 2022, Article 108539.
- [20] D.R. Brillinger, *Time Series: Data Analysis and Theory*, Expanded edition. New York: McGraw Hill, 1981.
- [21] 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.
- [22] N. Simon, J. Friedman, T. Hastie and R. Tibshirani, "A sparse-group lasso," *J. Computational Graphical Statistics*, vol. 22, pp. 231-245, 2013.
- [23] J. Gorski, F. Pfeuffer and K. Klamroth, "Biconvex sets and optimization with biconvex functions: A survey and extensions," *Mathematical Methods of Operations Research*, vol. 66, pp. 373-408, 2007.
- [24] 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.
- [25] 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.)
- [26] A.J. Rothman, P.J. Bickel, E. Levina and J. Zhu, "Sparse permutation invariant covariance estimation," *Electronic J. Statistics*, vol. 2, pp. 494-515, 2008.