

ESTIMATION OF DIFFERENTIAL GRAPHS FROM TIME-DEPENDENT DATA

Jitendra K. Tugnait

Department of Electrical & Computer Engineering
Auburn University, Auburn, AL 36849, USA. tugnait@auburn.edu

ABSTRACT

We consider the problem of estimating differences in two time series Gaussian graphical models (TSGGMs) which are known to have similar structure. The TSGGM structure is encoded in its inverse power spectral density (IPSD) just as the vector GGM structure is encoded in its precision (inverse covariance) matrix. Motivated by many applications, in existing works one is interested in estimating the difference in two precision matrices to characterize underlying changes in conditional dependencies of two sets of data comprised of independent and identically distributed observations. In this paper we consider estimation of the difference in two IPSD's to characterize underlying changes in conditional dependencies of two sets of time-dependent data. We analyze a group lasso penalized D-trace loss function approach in the frequency domain for differential graph learning, using Wirtinger calculus. An alternating direction method of multipliers (ADMM) algorithm is presented to optimize the objective function. Theoretical analysis establishing consistency of IPSD difference estimator in high-dimensional settings is presented. We illustrate our approach using a numerical example.

Keywords: Sparse graph learning; differential graph estimation; undirected graph; time series graphs.

1. INTRODUCTION

Graphical models are an important and useful tool for analyzing multivariate data [1]. A central concept is that of conditional independence. Given a collection of random variables, one wishes to assess the relationship between two variables, conditioned on the remaining variables. Consider a 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 $\mathcal{E} \subseteq [p] \times [p]$. Also consider a stationary (real-valued), zero-mean, p -dimensional multivariate Gaussian time series $\mathbf{x}(t)$, $t = 0, \pm 1, \pm 2, \dots$, with i th component $x_i(t)$, and correlation (covariance) matrix function $\mathbf{R}_{xx}(\tau) = \mathbb{E}\{\mathbf{x}(t + \tau)\mathbf{x}^T(t)\}$, $\tau = 0, \pm 1, \dots$. Given $\{\mathbf{x}(t)\}$, in the corresponding graph \mathcal{G} , each component series $\{x_i(t)\}$ is represented by a node (i in V), and associations between components $\{x_i(t)\}$ and $\{x_j(t)\}$ are represented by edges between nodes i and j of \mathcal{G} . In a conditional independence graph (CIG), there is no edge between nodes i and j if and only if (iff) $x_i(t)$ and $x_j(t)$ are conditionally independent given the remaining $p-2$ scalar series $x_\ell(t)$, $\ell \in [p]$, $\ell \neq i$, $\ell \neq j$.

A key insight in [2] was to transform the series to the frequency domain and express the graph relationships in the frequency domain. Denote the power spectral density (PSD) matrix of $\{\mathbf{x}(t)\}$ by $\mathbf{S}_x(f)$, where $\mathbf{S}_x(f) = \sum_{\tau=-\infty}^{\infty} \mathbf{R}_{xx}(\tau)e^{-j2\pi f\tau}$. In [2] 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 (IPSD), that is, $\{i, j\} \notin \mathcal{E}$ iff the (i, j) -th element of

$\mathbf{S}_x^{-1}(f)$ vanishes, i.e., $[\mathbf{S}_x^{-1}(f)]_{ij} = 0$ for every f . Hence one can use estimated IPSD of observed time series to infer the associated graph.

Graphical models were originally developed for random vectors (i.i.d. time series) [3, p. 234]. In particular, Gaussian graphical models (GGMs) are CIGs where \mathbf{x} is multivariate Gaussian. Suppose \mathbf{x} has positive-definite covariance matrix Σ with inverse covariance matrix $\Omega = \Sigma^{-1}$. Then Ω_{ij} , the (i, j) -th element of Ω , is zero iff x_i and x_j are conditionally independent. Such models have been extensively studied, and found to be useful in a wide variety of applications [4–7]. Graphical modeling of real-valued time-dependent data (stationary time series) originated with [8], followed by [2]. Nonparametric approaches for graphical modeling of real time series in high-dimensional settings (p is large and/or sample size n is of the order of p) have been investigated in [9–12], among others.

More recently there has been increasing interest in differential network analysis where one is interested in estimating the difference in two inverse covariance matrices [13–17]. Given observations \mathbf{x} and \mathbf{y} from two groups of subjects, one is interested in the difference $\Delta = \Omega_y - \Omega_x$, where $\Omega_x = (E\{\mathbf{x}\mathbf{x}^T\})^{-1}$ and $\Omega_y = (E\{\mathbf{y}\mathbf{y}^T\})^{-1}$. The associated differential graph is $\mathcal{G}_\Delta = (V, \mathcal{E}_\Delta)$ where $\{i, j\} \in \mathcal{E}_\Delta$ iff $[\Delta]_{ij} \neq 0$. It characterizes differences between the GGMs of the two sets of data. We use the term differential graph as in [17, 18] ([13, 14, 16] use the term differential network). As noted in [16], in biostatistics, the differential network/graph describes the changes in conditional dependencies between components under different environmental or genetic conditions. For instance, one may be interested in the differences in the graphical models of healthy and impaired subjects, or models under different disease states, given gene expression data or functional MRI signals [4, 19, 20].

In contrast to the approaches of [13–17], in this paper we address the problem of estimating differences in two time series Gaussian graphical models (TSGGMs) which are known to have similar structure. The TSGGM structure is encoded in its IPSD just as the vector GGM structure is encoded in its precision matrix. We consider estimation of the difference in two IPSD's to characterize underlying changes in conditional dependencies of two sets of time-dependent data $\{\mathbf{x}(t)\}_{t=1}^{n_x}$ and $\{\mathbf{y}(t)\}_{t=1}^{n_y}$. We analyze a group lasso penalized D-trace loss function approach in the frequency domain for differential graph learning, using Wirtinger calculus [21]. As a preliminary step, we first address the problem of estimation of complex differential graphs, given two complex-valued i.i.d. time series. This problem and the general problem of differential times series graph estimation have not been investigated before. The work of [18] considers time series differential graphs except that in [18] $\mathbf{x}(t)$ and $\mathbf{y}(t)$ are non-stationary (“functional” modeling), and instead of a single record (sample) of $\mathbf{x}(t)$, $t = 1, 2, \dots, n_x$ and $\mathbf{y}(t)$, $t = 1, 2, \dots, n_y$, as in this paper, they assume multiple independent observations of $\mathbf{x}(t)$, $t \in \mathcal{T}$, and $\mathbf{y}(t)$, $t \in \mathcal{T}$ (a closed subset of

This work is supported by NSF Grant CCF-2308473.

real line).

Notation: For a set V , $|V|$ denotes its cardinality. Given $\mathbf{A} \in \mathbb{C}^{p \times p}$, we use $\phi_{\min}(\mathbf{A})$, $\phi_{\max}(\mathbf{A})$, $|\mathbf{A}|$ and $\text{tr}(\mathbf{A})$ to denote the minimum eigenvalue, maximum eigenvalue, determinant and trace of \mathbf{A} , respectively. $[\mathbf{B}]_{ij}$ denotes the (i, j) -th element of \mathbf{B} , and so does B_{ij} . \mathbf{I} is the identity matrix. The symbol \otimes denotes the Kronecker product. The superscripts $*$ and H denote the complex conjugate and the Hermitian (conjugate transpose) operations, respectively. For $\mathbf{B} \in \mathbb{C}^{p \times q}$, we define $\|\mathbf{B}\| = \sqrt{\phi_{\max}(\mathbf{B}^H \mathbf{B})}$, $\|\mathbf{B}\|_F = \sqrt{\text{tr}(\mathbf{B}^H \mathbf{B})}$, $\|\mathbf{B}\|_1 = \sum_{i,j} |B_{ij}|$ and $\|\mathbf{B}\|_\infty = \max_{i,j} |B_{ij}|$. The notation $\mathbf{x} \sim \mathcal{N}_c(\mathbf{m}, \Sigma)$ denotes a random vector \mathbf{x} that is circularly symmetric (proper) complex Gaussian with mean \mathbf{m} and covariance Σ . Similarly, $\mathbf{x} \sim \mathcal{N}_r(\mathbf{m}, \Sigma)$ denotes a random vector \mathbf{x} that is real-valued Gaussian with mean \mathbf{m} and covariance Σ .

2. COMPLEX DIFFERENTIAL GRAPHS

We first recall a formulation of [13–16] for real-valued data. Let $\mathbf{x} \in \mathbb{R}^p$, $\mathbf{x} \sim \mathcal{N}_r(\mathbf{0}, \Sigma_{0x})$ and suppose we are given i.i.d. samples $\mathbf{x}(t)$, $t = 1, 2, \dots, n_x$, of \mathbf{x} , and similarly given i.i.d. samples $\mathbf{y}(t)$, $t = 1, 2, \dots, n_y$, of independent $\mathbf{y} \in \mathbb{R}^p$, $\mathbf{y} \sim \mathcal{N}_r(\mathbf{0}, \Sigma_{0y})$. Form the sample covariance estimates $\hat{\Sigma}_x = \frac{1}{n_x} \sum_{t=1}^{n_x} \mathbf{x}(t)\mathbf{x}^T(t)$ and $\hat{\Sigma}_y = \frac{1}{n_y} \sum_{t=1}^{n_y} \mathbf{y}(t)\mathbf{y}^T(t)$. In [13–16] one seeks to estimate $\Delta_0 = \Omega_{0y} - \Omega_{0x}$ and graph $\mathcal{G}_\Delta = (V, \mathcal{E}_\Delta)$, based on $\hat{\Sigma}_x$ and $\hat{\Sigma}_y$, where $\Omega_{0y} = \Sigma_{0y}^{-1}$, $\Omega_{0x} = \Sigma_{0x}^{-1}$. In [14] (see also [15, Sec. 2.1]), the following convex D-trace loss function is used

$$L_r(\Delta, \hat{\Sigma}_x, \hat{\Sigma}_y) = \frac{1}{2} \text{tr}(\hat{\Sigma}_x \Delta \hat{\Sigma}_y \Delta^T) - \text{tr}(\Delta(\hat{\Sigma}_x - \hat{\Sigma}_y)) \quad (1)$$

where D-trace refers to difference-in-trace loss function, a term coined in [22] in the context of graphical model estimation. The function $L_r(\Delta, \Sigma_{0x}, \Sigma_{0y})$ is strictly convex in Δ and has a unique minimum at $\Delta_0 = \Omega_{0y} - \Omega_{0x}$ [14, 15]. When one uses sample covariances, Δ is estimated by minimizing a lasso-penalized D-trace loss function [13–16] (group-lasso in [17]).

2.1. Proper Complex Gaussian Vectors

Consider complex-valued $\mathbf{x} \sim \mathcal{N}_c(\mathbf{0}, \Sigma_{0x})$ and $\mathbf{y} \sim \mathcal{N}_c(\mathbf{0}, \Sigma_{0y})$ with $\Sigma_{0x} \succ \mathbf{0}$, $\Sigma_{0y} \succ \mathbf{0}$. We need to estimate $\Delta_0 = \Omega_{0y} - \Omega_{0x}$. Consider the real-valued cost

$$L(\Delta, \hat{\Sigma}_x, \hat{\Sigma}_y) = \frac{1}{2} \left(\text{tr}(\hat{\Sigma}_x \Delta \hat{\Sigma}_y \Delta^H) + \text{tr}(\hat{\Sigma}_x^* \Delta^* \hat{\Sigma}_y^* \Delta^T) \right) - \text{tr} \left(\Delta(\hat{\Sigma}_x - \hat{\Sigma}_y) + \Delta^*(\hat{\Sigma}_x^* - \hat{\Sigma}_y^*) \right) \quad (2)$$

with $\hat{\Sigma}_x = \frac{1}{n_x} \sum_{t=1}^{n_x} \mathbf{x}(t)\mathbf{x}^H(t)$ and $\hat{\Sigma}_y = \frac{1}{n_y} \sum_{t=1}^{n_y} \mathbf{y}(t)\mathbf{y}^H(t)$. Using Wirtinger calculus, we find $\mathbf{0} = \frac{\partial L}{\partial \Delta^*} = \hat{\Sigma}_x \Delta \hat{\Sigma}_y - (\hat{\Sigma}_x - \hat{\Sigma}_y)$, implying $L(\Delta, \Sigma_{0x}, \Sigma_{0y})$ has a minimum at $\Delta_0 = \Omega_{0y} - \Omega_{0x}$. Define $\theta = [\text{vec}(\Delta)^T \text{vec}(\Delta)^H]^T \in \mathbb{C}^{2p^2}$. Then using $\text{tr}(\mathbf{A}^T \mathbf{B} \mathbf{C} \mathbf{D}^T) = \text{vec}(\mathbf{A})^T (\mathbf{D} \otimes \mathbf{B}) \text{vec}(\mathbf{C})$, we have $L(\Delta, \Sigma_{0x}, \Sigma_{0y}) = \frac{1}{2} \theta^H \mathcal{H} \theta - \theta^H \mathbf{b}$ where

$$\mathcal{H} = \begin{bmatrix} \Sigma_{0y}^* \otimes \Sigma_{0x} & \mathbf{0} \\ \mathbf{0} & \Sigma_{0y} \otimes \Sigma_{0x}^* \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \text{vec}(\Sigma_{0x} - \Sigma_{0y}) \\ \text{vec}(\Sigma_{0x}^* - \Sigma_{0y}^*) \end{bmatrix}.$$

As in the real case, $L(\Delta, \Sigma_{0x}, \Sigma_{0y})$ is strictly convex in Δ and the minimum at Δ_0 is unique since the Hessian $\mathcal{H} \succ \mathbf{0}$. For $\lambda > 0$,

define the lasso-penalized D-trace loss

$$L_\lambda(\Delta, \hat{\Sigma}_x, \hat{\Sigma}_y) = L(\Delta, \hat{\Sigma}_x, \hat{\Sigma}_y) + \lambda \sum_{k,\ell=1}^p |\Delta_{k\ell}|. \quad (3)$$

We seek $\hat{\Delta} = \arg \min_{\Delta} L_\lambda(\Delta, \hat{\Sigma}_x, \hat{\Sigma}_y)$.

2.2. Optimization

Similar to [15] (also [14]), we use an alternating direction method of multipliers (ADMM) approach [23] with variable splitting. Using variable splitting, consider

$$\min_{\Delta, \mathbf{W}} \left\{ L(\Delta, \hat{\Sigma}_x, \hat{\Sigma}_y) + \lambda \sum_{k,\ell=1}^p |W_{k\ell}| \right\} \text{ subject to } \Delta = \mathbf{W}. \quad (4)$$

The scaled augmented Lagrangian for this problem is [23]

$$L_\rho = L(\Delta, \hat{\Sigma}_x, \hat{\Sigma}_y) + \lambda \sum_{k,\ell=1}^p |W_{k\ell}| + \frac{\rho}{2} \|\Delta - \mathbf{W} + \mathbf{U}\|_F^2 \quad (5)$$

where \mathbf{U} is the dual variable, and $\rho > 0$ is the penalty parameter. Given the i th iteration results $\Delta^{(i)}$, $\mathbf{W}^{(i)}$, $\mathbf{U}^{(i)}$, in the $(i+1)$ st iteration, the algorithm executes the following 3 updates:

- (a) $\Delta^{(i+1)} \leftarrow \arg \min_{\Delta} L_a(\Delta)$, $L_a(\Delta) := L(\Delta, \hat{\Sigma}_x, \hat{\Sigma}_y) + \frac{\rho}{2} \|\Delta - \mathbf{W}^{(i)} + \mathbf{U}^{(i)}\|_F^2$
- (b) $\mathbf{W}^{(i+1)} \leftarrow \arg \min_{\mathbf{W}} L_b(\mathbf{W})$, $L_b(\mathbf{W}) := \lambda \sum_{k,\ell=1}^p |W_{k\ell}| + \frac{\rho}{2} \|\Delta^{(i+1)} - \mathbf{W} + \mathbf{U}^{(i)}\|_F^2$
- (c) $\mathbf{U}^{(i+1)} \leftarrow \mathbf{U}^{(i)} + \left(\Delta^{(i+1)} - \mathbf{W}^{(i+1)} \right)$

Update (a): Differentiate $L_a(\Delta)$ w.r.t. Δ^* to obtain $\mathbf{0} = \frac{\partial L_a(\Delta)}{\partial \Delta^*} = \hat{\Sigma}_x \Delta \hat{\Sigma}_y - (\hat{\Sigma}_x - \hat{\Sigma}_y) + \frac{\rho}{2} (\Delta - \mathbf{W}^{(i)} + \mathbf{U}^{(i)})$. Hence,

$$(\hat{\Sigma}_y^* \otimes \hat{\Sigma}_x + \frac{\rho}{2} \mathbf{I}) \text{vec}(\Delta) = \text{vec}(\hat{\Sigma}_x - \hat{\Sigma}_y + \frac{\rho}{2} (\mathbf{W}^{(i)} - \mathbf{U}^{(i)})) \quad (6)$$

where we used $\text{vec}(\mathbf{A} \mathbf{Y} \mathbf{B}) = (\mathbf{B}^T \otimes \mathbf{A}) \text{vec}(\mathbf{Y})$. Direct matrix inversion solution of (6) requires inversion of a $p^2 \times p^2$ matrix. A computationally cheaper solution follows as in [14, 15], but for Hermitian matrices. Carry out eigendecomposition of $\hat{\Sigma}_x$ and $\hat{\Sigma}_y$ as $\hat{\Sigma}_x = \mathbf{Q}_x \mathbf{D}_x \mathbf{Q}_x^H$, $\mathbf{Q}_x \mathbf{Q}_x^H = \mathbf{I}$ and $\hat{\Sigma}_y = \mathbf{Q}_y \mathbf{D}_y \mathbf{Q}_y^H$, $\mathbf{Q}_y \mathbf{Q}_y^H = \mathbf{I}$, where \mathbf{D}_x and \mathbf{D}_y are diagonal matrices. Then $\hat{\Delta}$ that minimizes $L_a(\Delta)$ is given by

$$\hat{\Delta} = \mathbf{Q}_x \left[\mathbf{D} \circ [\mathbf{Q}_x^H (\hat{\Sigma}_x - \hat{\Sigma}_y + \frac{\rho}{2} (\mathbf{W}^{(i)} - \mathbf{U}^{(i)})) \mathbf{Q}_y] \right] \mathbf{Q}_y^H \quad (7)$$

where the symbol \circ denotes the Hadamard product and $\mathbf{D} \in \mathbb{R}^{p \times p}$ organizes the diagonal of $(\mathbf{D}_y \otimes \mathbf{D}_x + \frac{\rho}{2} \mathbf{I})^{-1}$ in a matrix with $D_{jk} = 1 / ([\mathbf{D}_x]_{jj} [\mathbf{D}_y]_{kk} + \frac{\rho}{2})$. Note that the eigendecomposition of $\hat{\Sigma}_x$ and $\hat{\Sigma}_y$ has to be done only once. Thus

$$\Delta^{(i+1)} = \mathbf{Q}_x \left[\mathbf{D} \circ [\mathbf{Q}_x^H p(\hat{\Sigma}_x - \hat{\Sigma}_y + \frac{\rho}{2} (\mathbf{W}^{(i)} - \mathbf{U}^{(i)})) \mathbf{Q}_y] \right] \mathbf{Q}_y^H \quad (8)$$

Update (b): Here we have the lasso solution [12, Lemma 1]

$$W_{k\ell}^{(i+1)} = \left(1 - \frac{(\lambda/\rho)}{[|\Delta^{(i+1)} + \mathbf{U}^{(i)}|]_{k\ell}} \right)_+ [\Delta^{(i+1)} + \mathbf{U}^{(i)}]_{k\ell} \quad (9)$$

where $(a)_+ = \max(0, a)$. It results from separable optimization of $L_b(\mathbf{W}) = \sum_{k,\ell=1}^p \{\lambda |W_{k\ell}| + \frac{\rho}{2} |\Delta_{k\ell}^{(i+1)} - W_{k\ell} + U_{k\ell}^{(i)}|^2\}$.

Convergence. A stopping (convergence) criterion following [23, Sec. 3.3.1] can be devised. The stopping criterion is based on primal and dual residuals being small where, in our case, at $(i+1)$ st iteration, the primal residual is given by $\Delta^{(i+1)} - \mathbf{W}^{(i+1)}$ and the dual residual by $\rho(\mathbf{W}^{(i+1)} - \mathbf{W}^{(i)})$. Convergence criterion is met when the norms of these residuals are below some threshold. The objective function $L_\lambda(\Delta, \hat{\Sigma}_x, \hat{\Sigma}_y)$, given by (3), is strictly convex. It is also closed, proper and lower semi-continuous. Hence, for any fixed $\rho > 0$, the ADMM algorithm is guaranteed to converge [23, Sec. 3.2], in the sense that we have primal residual convergence to 0, dual residual convergence to 0, and objective function convergence to the optimal value.

3. DIFFERENTIAL TIME SERIES GRAPHS

We will address the problem via a frequency-domain formulation. For simplicity, we take $n_x = n_y = n$. Given $\mathbf{x}(t)$ and $\mathbf{y}(t)$ for $t = 1, 2, \dots, n$, define their respective (normalized) DFTs

$$\mathbf{d}_x(f_m) = \frac{1}{\sqrt{n}} \sum_{t=1}^n \mathbf{x}(t) \exp(-j2\pi f_m(t-1)), \quad (10)$$

$$\mathbf{d}_y(f_m) = \frac{1}{\sqrt{n}} \sum_{t=1}^n \mathbf{y}(t) \exp(-j2\pi f_m(t-1)), \quad (11)$$

$$f_m = \frac{m}{n}, \quad m = 0, 1, \dots, n-1. \quad (12)$$

The set of complex-valued random vectors $\{\mathbf{d}_x(f_m), \mathbf{d}_y(f_m)\}_{m=0}^{n/2}$ is a sufficient statistic for any inference problem based on data set $\{\mathbf{x}(t), \mathbf{y}(t)\}_{t=1}^n$ [12]. Suppose $\mathbf{S}_x(f_k)$ is locally smooth, so that $\mathbf{S}_x(f_k)$ is (approx.) constant over $K = 2m_t + 1$ consecutive frequency points f_m 's, $m_t > 0$. Pick $M = \lfloor (\frac{n}{2} - m_t - 1)/K \rfloor$ and $\tilde{f}_k = ((k-1)K + m_t + 1)/n$ for $k = 1, 2, \dots, M$ yielding M equally spaced frequencies \tilde{f}_k in the interval $(0, 0.5)$. It turns out that for "large" n , the DFT $\mathbf{d}_x(f_m)$ is a complex-valued proper (i.e., circularly symmetric) Gaussian vector $\sim \mathcal{N}_c(\mathbf{0}, \mathbf{S}_x(f_m))$, and (mutually) independent for $m = 1, 2, \dots, (n/2) - 1$, (n even) [24, Theorem 4.4.1], though not identically distributed, so long as the autocorrelation function of $\mathbf{x}(t)$ is absolutely summable. Similar comments apply to $\mathbf{d}_y(f_m)$. By local smoothness, $\mathbf{S}_x(\tilde{f}_{k,\ell}) = \mathbf{S}_x(\tilde{f}_k)$ for $\ell = -m_t, -m_t + 1, \dots, m_t$, where $\tilde{f}_{k,\ell} = \frac{(k-1)K + m_t + 1 + \ell}{n}$. Define

$$\hat{\mathbf{S}}_{xk} = \frac{1}{K} \sum_{\ell=-m_t}^{m_t} \mathbf{d}_x(\tilde{f}_{k,\ell}) \mathbf{d}_x^H(\tilde{f}_{k,\ell}) \quad (13)$$

$$\hat{\mathbf{S}}_{yk} = \frac{1}{K} \sum_{\ell=-m_t}^{m_t} \mathbf{d}_y(\tilde{f}_{k,\ell}) \mathbf{d}_y^H(\tilde{f}_{k,\ell}) \quad (14)$$

where $\hat{\mathbf{S}}_{xk}$ and $\hat{\mathbf{S}}_{yk}$ represent PSD estimators at frequency \tilde{f}_k using unweighted frequency-domain smoothing [24]. By local smoothness, $\mathbf{d}_x(\tilde{f}_{k,\ell}) \sim \mathcal{N}_c(\mathbf{0}, \mathbf{S}_x(\tilde{f}_k))$ and $\mathbf{d}_y(\tilde{f}_{k,\ell}) \sim \mathcal{N}_c(\mathbf{0}, \mathbf{S}_y(\tilde{f}_k))$.

Henceforth we will denote the true values of $\mathbf{S}_x(\tilde{f}_k)$ and $\mathbf{S}_y(\tilde{f}_k)$ as \mathbf{S}_{0xk} and \mathbf{S}_{0yk} , respectively, with their respective sample estimates $\hat{\mathbf{S}}_{xk}$ and $\hat{\mathbf{S}}_{yk}$, $k = 1, 2, \dots, M$. Let $\Delta_{0k} = \mathbf{S}_{0yk}^{-1} - \mathbf{S}_{0xk}^{-1}$. By local smoothness assumption, we apply the approach of Sec. 2.1 to estimate the differences Δ_{0k} , $k = 1, 2, \dots, M$, together with

a group-lasso penalty to enforce a common edgeset across the M frequency points. Define

$$\tilde{\Delta} = [\Delta_1, \Delta_2, \dots, \Delta_M] \in \mathbb{C}^{p^2 M}, \quad (15)$$

$$\Delta^{(ij)} = [[\Delta_1]_{ij}, [\Delta_2]_{ij}, \dots, [\Delta_M]_{ij}] \in \mathbb{C}^M. \quad (16)$$

We propose to estimate Δ_k 's by minimizing

$$L_f(\tilde{\Delta}) = \sum_{k=1}^M L(\Delta_k, \hat{\mathbf{S}}_{xk}, \hat{\mathbf{S}}_{yk}) + \lambda \sum_{i,j=1}^p \|\Delta^{(ij)}\| \quad (17)$$

In the estimated differential graph, $\{i, j\} \in \mathcal{E}_{\hat{\Delta}} \Leftrightarrow \|\hat{\Delta}^{(ij)}\| \neq 0$.

3.1. Optimization

As in Sec. 2.2, we use an ADMM approach. Using variable splitting, the scaled augmented Lagrangian for this problem is

$$L_\rho(\tilde{\Delta}, \tilde{\mathbf{W}}, \tilde{\mathbf{U}}) = L_f(\tilde{\Delta}) + \lambda \sum_{i,j=1}^p \|\mathbf{W}^{(ij)}\| + \frac{\rho}{2} \sum_{k=1}^M \|\Delta_k - \mathbf{W}_k + \mathbf{U}_k\|_F^2 \quad (18)$$

where $\tilde{\mathbf{W}} = [\mathbf{W}_1 \dots \mathbf{W}_M]$, $\tilde{\mathbf{U}} = [\mathbf{U}_1 \dots \mathbf{U}_M]$ is the dual variable, and $\rho > 0$ is the penalty parameter. Given the results $\tilde{\Delta}^{(m)}$, $\tilde{\mathbf{W}}^{(m)}$, $\tilde{\mathbf{U}}^{(m)}$ of the m th iteration, in the $(m+1)$ st iteration, an ADMM algorithm executes the following three updates:

- $\tilde{\Delta}^{(m+1)} \leftarrow \arg \min_{\tilde{\Delta}} \sum_{k=1}^M L_{ak}(\Delta_k)$, $L_{ak}(\Delta_k) = L(\Delta_k, \hat{\mathbf{S}}_{xk}, \hat{\mathbf{S}}_{yk}) + \frac{\rho}{2} \|\Delta_k - \mathbf{W}_k^{(m)} + \mathbf{U}_k^{(m)}\|_F^2$.
- $\tilde{\mathbf{W}}^{(m+1)} \leftarrow \arg \min_{\tilde{\mathbf{W}}} L_b(\tilde{\mathbf{W}})$, $L_b(\tilde{\mathbf{W}}) = \lambda \sum_{i,j=1}^p \|\mathbf{W}^{(ij)}\| + \frac{\rho}{2} \sum_{k=1}^M \|\Delta_k^{(m+1)} - \mathbf{W}_k^{(m)} + \mathbf{U}_k^{(m)}\|_F^2$.
- $\tilde{\mathbf{U}}^{(m+1)} \leftarrow \tilde{\mathbf{U}}^{(m)} + \tilde{\Delta}^{(m+1)} - \tilde{\mathbf{W}}^{(m+1)}$.

Update (a): Optimization in step (a) is separable in Δ_k , and the solution discussed in Sec. 2.2 applies. Therefore, with notational changes, for $k = 1, 2, \dots, M$, we have

$$\Delta_k^{(m+1)} = \mathbf{Q}_{xk} \left[\mathbf{D}^{(k)} \circ [\mathbf{Q}_{xk}^H (\hat{\mathbf{S}}_{xk} - \hat{\mathbf{S}}_{yk}) + \frac{\rho}{2} (\mathbf{W}_k^{(m)} - \mathbf{U}_k^{(m)})] \mathbf{Q}_{yk}^H \right] \mathbf{Q}_{yk} \quad (19)$$

where we have the eigendecomposition of $\hat{\mathbf{S}}_{xk}$ and $\hat{\mathbf{S}}_{yk}$ as $\hat{\mathbf{S}}_{xk} = \mathbf{Q}_{xk} \mathbf{D}_{xk} \mathbf{Q}_{xk}^H$ and $\hat{\mathbf{S}}_{yk} = \mathbf{Q}_{yk} \mathbf{D}_{yk} \mathbf{Q}_{yk}^H$, and $\mathbf{D}^{(k)} \in \mathbb{R}^{p \times p}$ organizes the diagonal of $(\mathbf{D}_{yk} \otimes \mathbf{D}_{xk} + \frac{\rho}{2} \mathbf{I})^{-1}$ in a matrix with $[\mathbf{D}^{(k)}]_{ij} = 1/([\mathbf{D}_{xk}]_{ii} [\mathbf{D}_{yk}]_{ii} + \frac{\rho}{2})$.

Update (b): Optimization in step (b) is separable in $\mathbf{W}^{(ij)}$, and the group lasso solution [12, Lemma 1] applies. With $\mathbf{A}_k = \Delta_k^{(m+1)} + \mathbf{U}_k^{(m)}$, and for $k = 1, \dots, M$ and $i, j = 1, \dots, p$,

$$[\mathbf{W}_k^{(m+1)}]_{ij} = \left(1 - \frac{\lambda}{\rho \|\mathbf{A}^{(ij)}\|} \right)_+ [\mathbf{A}_k]_{ij}, \quad (20)$$

$$\text{where } \mathbf{A}^{(ij)} = [[\mathbf{A}_1]_{ij}, \dots, [\mathbf{A}_M]_{ij}] \in \mathbb{C}^M. \quad (21)$$

4. THEORETICAL ANALYSIS

Here we analyze consistency of $\hat{\Delta}$ by following the approach of [26]. The i.i.d. results of [14, 15, 22] follow the method of [25] which

requires an irrepresentability condition that we do not impose. Define the true differential edgeset

$$\mathcal{E}_{\Delta_0} = \left\{ \{i, j\} : [S_{0y}^{-1}(f) - S_{0x}^{-1}(f)]_{ij} \neq 0, \right. \\ \left. i \neq j, 0 \leq f \leq 0.5 \right\}, \quad s = |\mathcal{E}_{\Delta_0}|. \quad (22)$$

Let $\mathbf{R}_{xx}(\tau) = \mathbb{E}\{\mathbf{x}(t + \tau)\mathbf{x}^T(t)\}$ and $\mathbf{R}_{yy}(\tau) = \mathbb{E}\{\mathbf{y}(t + \tau)\mathbf{y}^T(t)\}$. In the rest of this section, we allow $p, K = 2m_t + 1, M, s$ and λ to be a functions of sample size n , denoted as p_n, K_n, M_n, s_n and λ_n , respectively.

Define

$$M_0 = \max \left\{ \max_{f \in [0, 0.5]} \max_{ij} |[S_{0x}(f)]_{ij}|, \right. \\ \left. \max_{f \in [0, 0.5]} \max_{ij} |[S_{0y}(f)]_{ij}| \right\}, \quad (23)$$

$$\phi_{0, \min} = \min \left\{ \min_{f \in [0, 0.5]} \phi_{\min}(S_{0x}(f)), \min_{f \in [0, 0.5]} \phi_{\min}(S_{0y}(f)) \right\}, \quad (24)$$

$$M_d = \max_{f \in [0, 0.5]} \max_{ij} |[S_{0y}^{-1}(f) - S_{0x}^{-1}(f)]_{ij}|, \quad (25)$$

$$C_0 = 80 \max_{\ell, f} \left\{ [S_{0x}(f)]_{\ell\ell}, [S_{0y}(f)]_{\ell\ell} \right\} \\ \times \sqrt{2(\ln(16p_n^2 M_n) / \ln(p_n))} \quad (26)$$

where $\tau > 2$.

Let $\hat{\Delta} = \arg \min_{\tilde{\Delta}} L_f(\tilde{\Delta})$.

Theorem 1 : Assume that $\sum_{\tau=-\infty}^{\infty} \|\mathbf{R}_{xx}(\tau)\|_{k\ell} < \infty$ and $\sum_{\tau=-\infty}^{\infty} \|\mathbf{R}_{yy}(\tau)\|_{k\ell} < \infty$ for every $k, \ell \in [p]$, Under (22), if

$$\lambda_n \geq 2\sqrt{M_n}(3M_0 s_n M_d + 2)C_0 \sqrt{\frac{\ln(p_n)}{K_n}}, \quad (27)$$

$$K_n > \left(\frac{96M_n M_0 s_n}{\phi_{0, \min}} \right)^2 C_0^2 \ln(p_n), \quad (28)$$

then with probability $> 1 - 2/p_n^{\tau-2}$, for any $\tau > 2$, we have

$$\|\hat{\Delta} - \tilde{\Delta}_0\|_F \leq \frac{12\sqrt{s_n} \lambda_n}{\phi_{0, \min}} \bullet \quad (29)$$

The proof of Theorem 1 is omitted for lack of space.

Remark 1: Convergence Rate. If $M_0, M_d, \phi_{0, \min}$ and C_0 stay bounded with increasing sample size n , we have $\|\hat{\Delta} - \tilde{\Delta}_0\|_F = \mathcal{O}_P(s_n^{1.5} \sqrt{M_n \ln(p_n) / K_n})$. Therefore, for $\|\hat{\Delta} - \tilde{\Delta}_0\|_F \rightarrow 0$ as $n \rightarrow \infty$, we must have $s_n^{1.5} \sqrt{M_n \ln(p_n) / K_n} \rightarrow 0$. \square

5. NUMERICAL EXAMPLE

Consider $p = 128$, 16 clusters (communities) of 8 nodes each, where nodes within a community are not connected to any nodes in other communities. Within any community of 8 nodes, the x -data are generated using a vector autoregressive (VAR) model of order 3. Consider community $q, q = 1, 2, \dots, 16$. Then $\mathbf{x}^{(q)}(t) \in \mathbb{R}^8$ is generated as $\mathbf{x}^{(q)}(t) = \sum_{i=1}^3 \mathbf{A}_i^{(q)} \mathbf{x}^{(q)}(t-i) + \mathbf{w}^{(q)}(t)$. Only 15% of entries of $\mathbf{A}_i^{(q)}$'s are nonzero and the nonzero elements are independently and uniformly distributed over $[-0.6, 0.6]$. 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 overall data $\mathbf{x}(t)$ is given by $\mathbf{x}(t) = [\mathbf{x}^{(1)\top}(t) \dots \mathbf{x}^{(16)\top}(t)]^\top \in \mathbb{R}^p$ with $\mathbf{w}^{(q)}(t)$ as i.i.d. zero-mean Gaussian with identity covariance matrix. To generate y -data, we randomly eliminate one of the 16 clusters of $\mathbf{x}(t)$ and replace it with an independently generated $\mathbf{y}^{(q)}(t)$ mimicking generation of $\mathbf{x}^{(q)}(t)$. First 100 samples are discarded to eliminate transients. This set-up leads to a differential time series graph with one cluster difference (8^2 edges out of 128^2 edges). We generate $n = n_x = n_y$ observations for $\mathbf{x}(t)$ and $\mathbf{y}(t)$, with $n \in \{512, 2048, 4096\}$.

Simulation results based on 100 runs are shown in Fig. 1. By changing the penalty parameter λ and determining the resulting edges over 100 runs, we calculated the true positive rate (TPR) which calculates true edges correctly detected ($\|\hat{\Delta}^{(ij)}\| \neq 0$ and $\|\tilde{\Delta}_0^{(ij)}\| \neq 0$), and false positive rate 1-TNR (where TNR is the true negative rate) which are the edges $\{i, j\}$ for which $\|\hat{\Delta}^{(ij)}\| \neq 0$ but $\|\tilde{\Delta}_0^{(ij)}\| = 0$. (Estimated $\hat{\Delta}_k$'s are not necessarily Hermitian. We use $(\hat{\Delta}_k + \hat{\Delta}_k^H)/2$ as the Hermitian estimate.) The receiver operating characteristic (ROC) is shown in Fig. 1 for our proposed approach (labeled "DTS") as well as for an approach that assumes the data is i.i.d. (labeled "IID"), based on [15] which minimizes lasso-penalized (1) based on difference of precision matrices. For the proposed approach we used $M = 2$ and $K = 127, 511, 1023$ for $n = 512, 2048, 4096$, respectively. It is seen from Fig. 1 that our approach significantly outperforms the IID approach, yielding much higher TPR for a given 1-TNR.

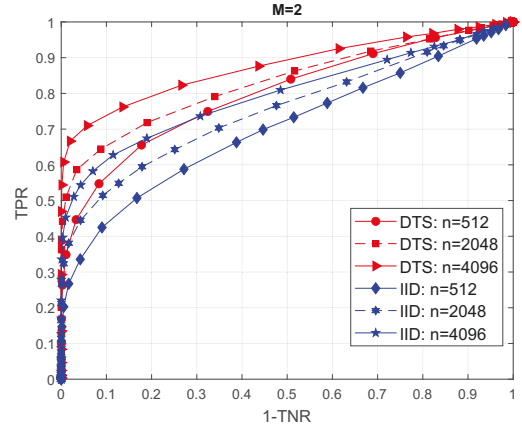


Fig. 1: ROC curves: DTS is the proposed approach and IID is the approach of [15]. TPR=true positive rate, TNR=true negative rate

6. CONCLUSIONS

We addressed the problem of estimating differences in two time series Gaussian graphical models (TSGGMs) which are known to have similar structure, via estimation of the difference in their IPSD's. We analyzed a group lasso penalized D-trace loss function approach in the frequency domain. An ADMM algorithm was presented to optimize the convex objective function. Theoretical analysis establishing consistency of the estimator of IPSD difference in high-dimensional settings was performed. We illustrated our approach via a numerical example.

7. REFERENCES

- [1] S.L. Lauritzen, *Graphical models*. Oxford, UK: Oxford Univ. Press, 1996.

- [2] R. Dahlhaus, "Graphical interaction models for multivariate time series," *Metrika*, vol. 51, pp. 157-172, 2000.
- [3] M. Eichler, "Graphical modelling of multivariate time series," *Probability Theory and Related Fields*, vol. 153, issue 1-2, pp. 233-268, June 2012.
- [4] 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.
- [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] N. Meinshausen and P. Bühlmann, "High-dimensional graphs and variable selection with the Lasso," *Ann. Statist.*, vol. 34, no. 3, pp. 1436-1462, 2006.
- [7] 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.
- [8] D.R. Brillinger, "Remarks concerning graphical models of times series and point processes," *Revista de Econometria (Brazilian Rev. Econometr.)*, vol. 16, pp. 1-23, 1996.
- [9] A. Jung, "Learning the conditional independence structure of stationary time series: A multitask learning approach," *IEEE Trans. Signal Process.*, vol. 63, no. 21, pp. 5677-5690, Nov. 1, 2015.
- [10] 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.
- [11] J.K. Tugnait, "Consistency of sparse-group lasso graphical model selection for time series," in *Proc. 54th Asilomar Conference on Signals, Systems and Computers*, Pacific Grove, CA, Nov. 1-4, 2020, pp. 589-593.
- [12] 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.
- [13] Y. Wu, T. Li, X. Liu and L.I. Chen, "Differential network inference via the fused D-trace loss with cross variables," *Electronic J. Statistics*, vol. 14, pp. 1269-1301, 2020.
- [14] H. Yuan, R. Xi, C. Chen and M. Deng, "Differential network analysis via lasso penalized D-trace loss," *Biometrika*, vol. 104, pp. 755-770, 2017.
- [15] B. Jiang, X. Wang and C. Leng, "A direct approach for sparse quadratic discriminant analysis," *J. Machine Learning Research*, vol. 19, pp. 1-37, 2018.
- [16] Z. Tang, Z. Yu and C. Wang, "A fast iterative algorithm for high-dimensional differential network," *Computational Statistics*, vol. 35, pp. 95-109, 2020.
- [17] J.K. Tugnait, "Estimation of high-dimensional differential graphs from multi-attribute data," in *Proc. 2023 IEEE Intern. Conf. Acoustics, Speech & Signal Processing (ICASSP 2023)*, pp. 1-5, Rhodes Island, Greece, June 4-9, 2023.
- [18] B. Zhao, Y.S. Wang and M. Kolar, "FuDGE: A method to estimate a functional differential graph in a high-dimensional setting," *J. Machine Learning Research*, vol. 23, pp. 1-82, 2022.
- [19] S.D. Zhao, T.T. Cai and H. Li, "Direct estimation of differential networks," *Biometrika*, vol. 101, pp. 253-268, June 2014.
- [20] E. Belilovsky, G. Varoquaux and M.B. Blaschko, "Hypothesis testing for differences in Gaussian graphical models: Applications to brain connectivity," *Advances in Neural Information Processing Systems (NIPS 2016)*, vol. 29, Dec. 2016.
- [21] P.J. Schreier and L.L. Scharf, *Statistical Signal Processing of Complex-Valued Data*, Cambridge, UK: Cambridge Univ. Press, 2010.
- [22] T. Zhang and H. Zou, "Sparse precision matrix estimation via lasso penalized D-trace loss," *Biometrika*, vol. 101, pp. 103-120, 2014.
- [23] 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.
- [24] D.R. Brillinger, *Time Series: Data Analysis and Theory*, Expanded edition. New York: McGraw Hill, 1981.
- [25] P. Ravikumar, M.J. Wainwright, G. Raskutti and B. Yu, "High-dimensional covariance estimation by minimizing ℓ_1 -penalized log-determinant divergence," *Electronic J. Statistics*, vol. 5, pp. 935-980, 2011.
- [26] S.N. Negahban, P. Ravikumar, M.J. Wainwright and B. Yu, "A unified framework for high-dimensional analysis of M-estimators with decomposable regularizers," *Statistical Science*, vol. 27, No. 4, pp. 538-557, 2012.