

Consistency of Sparse-Group Lasso Graphical Model Selection for Time Series

Jitendra K. Tugnait

Dept. of Electrical & Computer Eng.
Auburn University, Auburn, AL 36849, USA

Abstract—We consider the problem of inferring the conditional independence graph (CIG) of a high-dimensional stationary multivariate Gaussian time series. A p -variate Gaussian time series graphical model associated with an undirected graph with p vertices is defined as the family of time series that obey the conditional independence restrictions implied by the edge set of the graph. A sparse-group lasso-based frequency-domain formulation of the problem has been considered in the literature where the objective is to estimate the inverse power spectral density (PSD) of the data via optimization of a sparse-group lasso based penalized log-likelihood cost function that is formulated in the frequency-domain. The CIG is then inferred from the estimated inverse PSD. In this paper we establish sufficient conditions for consistency of the inverse PSD estimator resulting from the sparse-group graphical lasso-based approach.

I. INTRODUCTION

Graphical interaction models (“graphical models,” in short) are an important and useful tool for analyzing multivariate data [1]. Graphical modeling is a form of multivariate analysis where one uses graphs to represent models. 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. In graphical models, graphs are used to display the conditional independence structure of the 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 $x_\ell(t)$, $\ell \in [p]$, $\ell \neq i$, $\ell \neq j$. Thus, edge $\{i, j\} \in \mathcal{E}$ iff time series components $x_i(t)$ and $x_j(t)$ are conditionally dependent, and edge $\{i, j\} \notin \mathcal{E}$ iff $x_i(t)$ and $x_j(t)$ are conditionally independent. Gaussian graphical models (GGM) are CIGs where $\{\mathbf{x}(t)\}$ is a multivariate Gaussian sequence.

A key insight in [2], [3] 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 $\mathbf{S}_x(f) = \sum_{\tau=-\infty}^{\infty} \mathbf{R}_{xx}(\tau) e^{-j2\pi f\tau}$, the Fourier transform of $\mathbf{R}_{xx}(\tau)$. Here f is the normalized frequency, in Hz, in the interval $[0, 1)$ or $(-0.5, 0.5]$. In [2], [3] 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 .

Graphical models were originally developed for random vectors with multiple independent realizations [4, p. 234], i.e., for time series that is independent and identically distributed (i.i.d.): p -dimensional $\mathbf{x}(t)$, $t = 1, 2, \dots$, with $\mathbf{x}(t_1)$ independent of $\mathbf{x}(t_2)$ for $t_1 \neq t_2$, and $\mathbf{x}(t)$ identically distributed for any (integer) t . Such models have been extensively studied, and found to be useful in a wide variety of applications [5]–[10]. Graphical modeling of real-valued time-dependent data (stationary time series) originated with [2], followed by [3].

A sparse-group lasso-based frequency-domain formulation of the problem was investigated in [11] where the objective was to estimate the inverse power spectral density (PSD) of the data via optimization of a sparse-group lasso based penalized log-likelihood cost function that was formulated in the frequency-domain. The CIG is then inferred from the estimated inverse PSD. In this paper we establish sufficient conditions for consistency of the approach of [11]. Only the computational aspects of this problem were addressed in [11] where simulation comparisons with [12] were also provided; [11] significantly outperforms [12]. Further comparisons with [12] are in [11].

Notation: Given $\mathbf{A} \in \mathbb{C}^{p \times p}$, we use $\phi_{\min}(\mathbf{A})$, $\phi_{\max}(\mathbf{A})$, $|\mathbf{A}|$, $\text{tr}(\mathbf{A})$ and $\text{etr}(\mathbf{A})$ to denote the minimum eigenvalue, maximum eigenvalue, determinant, trace, and exponential of trace of \mathbf{A} , respectively. The Kronecker product of matrices \mathbf{A} and \mathbf{B} is denoted by $\mathbf{A} \otimes \mathbf{B}$. For $\mathbf{B} \in \mathbb{C}^{p \times q}$, we define $\|\mathbf{B}\| = \sqrt{\phi_{\max}(\mathbf{B}^H \mathbf{B})}$, $\|\mathbf{B}\|_F = \sqrt{\text{tr}(\mathbf{B}^H \mathbf{B})}$ and $\|\mathbf{B}\|_1 = \sum_{i,j} |B_{ij}|$ where B_{ij} is the (i, j) -th element of \mathbf{B} , also denoted by $[\mathbf{B}]_{ij}$. 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 $(\mathbf{A}^\top)^{-1}$. 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 $0 < M < \infty$ such that $P(\|\mathbf{y}_n\| \leq M\|\mathbf{x}_n\|) \geq 1 - \varepsilon \forall n \geq 1$.

II. PENALIZED LOG-LIKELIHOOD

Given $\mathbf{x}(t)$ for $t = 0, 1, 2, \dots, n-1$. Define the (normalized) DFT $\mathbf{d}_x(f_m)$ of $\mathbf{x}(t)$, ($j = \sqrt{-1}$),

$$\mathbf{d}_x(f_m) = \frac{1}{\sqrt{n}} \sum_{t=0}^{n-1} \mathbf{x}(t) \exp(-j2\pi f_m t) \quad (1)$$

where

$$f_m = m/n, \quad m = 0, 1, \dots, n-1. \quad (2)$$

It is established in [13] that the set of complex-valued random vectors $\{\mathbf{d}_x(f_m)\}_{m=0}^{n/2}$ is a sufficient statistic for any inference problem based on dataset $\{\mathbf{x}(t)\}_{t=0}^{n-1}$. Suppose $\mathbf{S}_x(f_m)$ is locally smooth (a standard

This work was supported by the National Science Foundation under Grants CCF-1617610 and ECCS-2040536. E-mail: tugnajk@auburn.edu

assumption in PSD estimation), so that $\mathbf{S}_x(f_m)$ is (approximately) constant over $K = 2m_t + 1$ consecutive frequency points f_m 's. Pick

$$\tilde{f}_k = \frac{(k-1)K + m_t + 1}{N}, \quad k = 1, 2, \dots, M, \quad (3)$$

$$M = \left\lfloor \frac{\frac{N}{2} - m_t - 1}{K} \right\rfloor, \quad (4)$$

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

$$\mathbf{S}_x(\tilde{f}_{k,\ell}) = \mathbf{S}_x(\tilde{f}_k) \text{ for } \ell = -m_t, -m_t + 1, \dots, m_t, \quad (5)$$

where

$$\tilde{f}_{k,\ell} = \frac{(k-1)K + m_t + 1 + \ell}{n}. \quad (6)$$

It is known ([14, Theorem 4.4.1]) that asymptotically (as $n \rightarrow \infty$), $\mathbf{d}_x(f_m)$, $m = 1, 2, \dots, (n/2) - 1$, (N even), are independent proper (i.e., circularly symmetric), complex Gaussian $\mathcal{N}_c(\mathbf{0}, \mathbf{S}_x(f_m))$ random vectors, respectively; $\mathbf{x}(t)$ need not be Gaussian but must satisfy some regularity conditions [13]. Then the joint probability density of the sufficient statistic, for large n , is

$$f_{\mathbf{D}}(\mathbf{D}) = \prod_{k=1}^M \frac{\exp\left(-\text{tr}(\tilde{\mathbf{D}}(\tilde{f}_k)\mathbf{S}_x^{-1}(\tilde{f}_k))\right)}{\pi^{Kp} |\mathbf{S}_x(\tilde{f}_k)|^K} = \prod_{k=1}^M f_{\tilde{\mathbf{D}}(\tilde{f}_k)}(\tilde{\mathbf{D}}(\tilde{f}_k)) \quad (7)$$

where

$$\tilde{\mathbf{D}}(\tilde{f}_k) = \left[\mathbf{d}_x(\tilde{f}_{k,-m_t}) \mathbf{d}_x(\tilde{f}_{k,-m_t+1}) \cdots \mathbf{d}_x(\tilde{f}_{k,m_t}) \right]^H \quad (8)$$

$$\hat{\mathbf{S}}_k = \frac{1}{K} \underbrace{\sum_{\ell=-m_t}^{m_t} \mathbf{d}_x(\tilde{f}_{k,\ell}) \mathbf{d}_x^H(\tilde{f}_{k,\ell})}_{=: \tilde{\mathbf{D}}(\tilde{f}_k)} \quad (9)$$

and $\hat{\mathbf{S}}_k$ represents PSD estimator at frequency \tilde{f}_k using unweighted frequency-domain smoothing.

In the high-dimension case of $K < p(p-1)/2$ (# of unknowns in $\mathbf{S}_x^{-1}(\tilde{f}_k)$), one may need to use penalty terms to enforce sparsity and to make the problem well-conditioned. We wish to estimate inverse PSD matrix $\Phi_k := \mathbf{S}_x^{-1}(\tilde{f}_k)$. In terms of Φ_k , we have the log-likelihood [11]

$$\ln f_{\mathbf{D}}(\mathbf{D}) \propto -G(\{\Phi\}) \quad (10)$$

$$:= \sum_{k=1}^M \left[(\ln |\Phi_k| + \ln |\Phi_k^*|) - \text{tr} \left(\hat{\mathbf{S}}_k \Phi_k + \hat{\mathbf{S}}_k^* \Phi_k^* \right) \right] \quad (11)$$

where the first expression in (11) follows by specifying the pdf of \mathbf{D} in terms of joint pdf of \mathbf{D} and \mathbf{D}^* (correct way to handle complex variates [15]). Imposing a sparse-group sparsity constraint [5], [16], minimize a penalized version of negative log-likelihood w.r.t. $\{\Phi\}$

$$L_{SGL}(\{\Phi\}) = -\ln f_{\mathbf{D}}(\mathbf{D}) + P(\{\Phi\}), \quad (12)$$

$$P(\{\Phi\}) := \bar{\lambda}_1 \sum_{k=1}^M \sum_{\substack{i,j=1 \\ i \neq j}}^p |\Phi_k]_{ij}| + \bar{\lambda}_2 \sum_{\substack{i,j=1 \\ i \neq j}}^p \sqrt{\sum_{k=1}^M |\Phi_k]_{ij}|^2} \quad (13)$$

where $\bar{\lambda}_1, \bar{\lambda}_2 \geq 0$ are tuning parameters. Computational aspects of this problem were addressed in [11] where simulation comparisons with [12] were also provided; [11] significantly outperforms [12].

In this paper we analyze properties of the minimizer $\{\hat{\Phi}\}$ of $L_{GL}(\{\Phi\})$.

III. CONSISTENCY

Define $p \times (pM)$ matrix Ω :

$$\Omega = [\Phi_1 \Phi_2 \cdots \Phi_M] \quad (14)$$

With $0 \leq \alpha \leq 1$, re-express the objective function (12) as

$$L_{SGL}(\Omega) = G(\{\Phi\}) + \alpha \lambda_n \sum_{k=1}^{M_n} \sum_{\substack{i,j=1 \\ i \neq j}}^{p_n} |\Phi_k]_{ij}| \\ + (1-\alpha) \lambda_n \sum_{\substack{i,j=1 \\ i \neq j}}^{p_n} \sqrt{\sum_{k=1}^{M_n} |\Phi_k]_{ij}|^2} \quad (15)$$

where we now allow p, M, K (see (3)), and λ to be functions of sample size n , denoted as p_n, M_n, K_n and λ_n , respectively. Note that $K_n M_n \approx n/2$. Pick $K_n = a_1 n^\gamma$ and $M_n = a_2 n^{1-\gamma}$ for some $0.5 < \gamma < 1$ so that $M_n/K_n \rightarrow 0$ as $n \rightarrow \infty$. We have rewritten $\bar{\lambda}_1$ and $\bar{\lambda}_2$ as $\alpha \lambda_n$ and $(1-\alpha) \lambda_n$, respectively, following [16]. Parameter $\lambda > 0$ is a penalty (tuning) parameter used to control sparsity, and $0 \leq \alpha \leq 1$ yields a convex combination of lasso and group lasso penalties ($\alpha = 0$ gives the group-lasso fit while $\alpha = 1$ yields the lasso fit). In (15), an ℓ_1 penalty term is applied to each off-diagonal element of Φ_k via $\alpha \lambda_n |\Phi_k]_{ij}|$ (lasso), and to the off-block-diagonal group of M_n terms via $(1-\alpha) \lambda_n \sqrt{\sum_{k=1}^{M_n} |\Phi_k]_{ij}|^2}$ (group lasso). The function $L_{SGL}(\Omega)$ is strictly convex in Ω for $\Phi_k \succ \mathbf{0} \forall k$.

We follow proof technique of [17] which deals with i.i.d time series models and lasso penalty, to establish our main result, Theorem 1. Assume

- (A1) Define the true edge set of the graph by \mathcal{E}_0 , implying that $\mathcal{E}_0 = \{\{i, j\} : [\mathbf{S}_0^{-1}(f)]_{ij} \neq 0, i \neq j, 0 \leq f \leq 0.5\}$ where $\mathbf{S}_0(f)$ denotes the true PSD of $\mathbf{x}(t)$. (We also use Φ_{0k} for $\mathbf{S}_0^{-1}(\tilde{f}_k)$ where \tilde{f}_k is as in (3), and use Ω_0 to denote the true value of Ω). Assume that $\text{card}(\mathcal{E}_0) = |\mathcal{E}_0| \leq s_{n0}$.
- (A2) The minimum and maximum eigenvalues of $p_n \times p_n$ PSD $\mathbf{S}_0(f) \succ \mathbf{0}$ satisfy

$$0 < \beta_{\min} \leq \min_{f \in [0, 0.5]} \phi_{\min}(\mathbf{S}_0(f)) \\ \leq \max_{f \in [0, 0.5]} \phi_{\max}(\mathbf{S}_0(f)) \leq \beta_{\max} < \infty.$$

Here β_{\min} and β_{\max} are not functions of n (or p_n).

Let $\hat{\Omega}_\lambda = \arg \min_{\Omega: \Phi_k \succ \mathbf{0}} L_{SGL}(\Omega)$. Theorem 1 whose proof is given in Sec. IV, establishes consistency of $\hat{\Omega}_\lambda$.

Theorem 1 (Consistency). For $\tau > 2$, let

$$C_0 = 80 \max_{\ell, f} ([\mathbf{S}_0(f)]_{\ell\ell}) \sqrt{\frac{2 \ln(16 p_n^\tau M_n)}{\ln(p_n)}}. \quad (16)$$

Given real numbers $\delta_1 \in (0, 1)$, $\delta_2 > 0$ and $C_1 > 0$, let

$$R = C_2 C_0 / \beta_{\min}^2, \quad C_2 = (4 + C_1 + \delta_2)(1 + \delta_1)^2, \quad (17)$$

$$r_n = \sqrt{\frac{M_n(p_n + s_{n0}) \ln(p_n)}{K_n}}, \quad C_2 r_n = o(1), \quad (18)$$

$$N_1 = 2 \ln(16 p_n^\tau M_n), \quad (19)$$

$$N_2 = \arg \min \left\{ n : r_n \leq \frac{\delta_1 \beta_{\min}}{C_2 C_0} \right\}. \quad (20)$$

Suppose the regularization parameter λ_n and $\alpha \in [0, 1]$ satisfy

$$\begin{aligned} 2C_0 \sqrt{\frac{\ln(p_n)}{K_n}} &\leq \frac{\lambda_n}{\sqrt{M_n}} \\ &\leq \frac{C_1 C_0}{1 + \alpha(\sqrt{M_n} - 1)} \sqrt{\left(1 + \frac{p_n}{s_{n0}}\right) \frac{\ln(p_n)}{K_n}}. \end{aligned} \quad (21)$$

Then if the sample size is such that $K_n > \max\{N_1, N_2\}$ and assumptions (A1)-(A2) hold true, $\hat{\Omega}_\lambda$ satisfies

$$\|\hat{\Omega}_\lambda - \Omega_0\|_F \leq Rr_n \quad (22)$$

with probability greater than $1 - 1/p_n^{\tau-2}$. A sufficient condition for the lower bound in (21) to be less than the upper bound for every $\alpha \in [0, 1]$ is $C_1 = 2(1 + \alpha(\sqrt{M_n} - 1))$. In terms of rate of convergence, $\|\hat{\Omega}_\lambda - \Omega_0\|_F = \mathcal{O}_P(C_1 r_n)$. •

Remark. If $\alpha = 0$, then C_1 is a constant, and therefore, $\|\hat{\Omega}_\lambda - \Omega_0\|_F = \mathcal{O}_P(r_n)$. If $\alpha > 0$, then $C_1 = \mathcal{O}(\sqrt{M_n})$, therefore, $\|\hat{\Omega}_\lambda - \Omega_0\|_F = \mathcal{O}_P(\sqrt{M_n} r_n)$. As note before, since $K_n M_n \approx n/2$, if one picks $K_n = a_1 n^\gamma$ and $M_n = a_2 n^{1-\gamma}$ for some $0 < \gamma < 1$, then we must have $\frac{2}{3} < \gamma < 1$ so that $M_n^2/K_n \rightarrow 0$ as $n \rightarrow \infty$, ensuring $C_2 r_n = o(1)$. In fact, one needs $\sqrt{\frac{M_n^2(p_n + s_{n0}) \ln(p_n)}{K_n}} \rightarrow 0$. □

IV. PROOF OF THEOREM 1

Our proof relies on the method of [17] which deals with i.i.d time series models and lasso penalty, and our prior results in [18] dealing with complex Gaussian vectors (not time series). From now on we use the term “with high probability” (w.h.p.) to denote with probability greater than $1 - 1/p_n^{\tau-2}$. First we need several auxiliary results.

Lemma 1 below is specialization of [19, Lemma 1] to Gaussian random vectors. It follows from [19, Lemma 1] after setting the sub-Gaussian parameter σ in [19, Lemma 1] to 1.

Lemma 1. Consider a zero-mean Gaussian random vector $\mathbf{z} \in \mathbb{R}^p$ with covariance $\mathbf{R} \succ \mathbf{0}$. Given n i.i.d. samples $\mathbf{z}(t)$, $t = 0, 1, \dots, n-1$, of \mathbf{z} , let $\hat{\mathbf{R}} = (1/n) \sum_{t=0}^{n-1} \mathbf{z}\mathbf{z}^\top$ denote the sample covariance matrix. Then $\hat{\mathbf{R}}$ satisfies the tail bound

$$P\left(\left|[\hat{\mathbf{R}} - \mathbf{R}]_{ij}\right| > \delta\right) \leq 4 \exp\left(-\frac{n\delta^2}{3200 \max_i(R_{ii}^2)}\right) \quad (23)$$

for all $\delta \in (0, 40 \max_i(R_{ii}))$. •

Exploiting Lemma 1, we have Lemma 2 regarding $\hat{\mathbf{S}}_k$. We denote $\mathbf{S}_0(f_k)$ as \mathbf{S}_{0k} in this section.

Lemma 2. Under Assumption (A2), $\hat{\mathbf{S}}_k$ satisfies the tail bound

$$P\left(\max_{k,q,l} \left|[\hat{\mathbf{S}}_k - \mathbf{S}_{0k}]_{ql}\right| > C_0 \sqrt{\frac{\ln(p_n)}{K_n}}\right) \leq \frac{1}{p_n^{\tau-2}} \quad (24)$$

for $\tau > 2$, if the sample size $n > N_1$, where C_0 is defined in (16) and N_1 is defined in (19). •

Proof. We will use $\mathbf{d}_k(\ell)$ for $\mathbf{d}_x(\tilde{f}_{k,\ell})$. Notice that $\hat{\mathbf{S}}_k$ is the covariance estimate based on K_n i.i.d. samples $\mathbf{d}_k(\ell)$ of $\mathbf{d}_k \sim \mathcal{N}_c(\mathbf{0}, \mathbf{S}_{0k})$. Define $\mathbf{d}_{kr} = \text{Re}(\mathbf{d}_k)$, $\mathbf{d}_{ki} = \text{Im}(\mathbf{d}_k)$ and $\mathbf{z}_k = [\mathbf{d}_{kr}^\top \ \mathbf{d}_{ki}^\top]^\top \in \mathbb{R}^{2p_n}$. Then with $\mathbf{R}_{y_1 y_2} := \mathbb{E}\{\mathbf{y}_1 \mathbf{y}_2^\top\}$ and $\hat{\mathbf{R}}_{y_1 y_2} := (1/K_n) \sum_\ell \mathbf{y}_1(\ell) \mathbf{y}_2^\top(\ell)$, we have

$$\hat{\mathbf{S}}_k = \hat{\mathbf{R}}_{d_{kr} d_{kr}} + \hat{\mathbf{R}}_{d_{ki} d_{ki}} + j \left(-\hat{\mathbf{R}}_{d_{kr} d_{ki}} + \hat{\mathbf{R}}_{d_{ki} d_{kr}} \right) \quad (25)$$

and

$$\mathbf{S}_{0k} = \mathbf{R}_{d_{kr} d_{kr}} + \mathbf{R}_{d_{ki} d_{ki}} + j \left(-\mathbf{R}_{d_{kr} d_{ki}} + \mathbf{R}_{d_{ki} d_{kr}} \right) \quad (26)$$

Also, $\mathbf{z}_k \sim \mathcal{N}_r(\mathbf{0}, \mathbf{R}_{0k})$. Since $\mathbf{d}_k(\ell)$ is proper, a little algebra leads to

$$[\mathbf{S}_{0k}]_{qq} = 2[\mathbf{R}_{0k}]_{qq}, \quad q = 1, 2, \dots, p_n \quad (27)$$

Denote the estimate of \mathbf{R}_{0k} based on K_n samples as $\hat{\mathbf{R}}_{zzk}$. By Lemma 1 and applying the union bound over all $M_n(2p_n)^2$ entries of $\hat{\mathbf{R}}_{zzk} - \mathbf{R}_{0k}$, $k = 1, 2, \dots, K_n$, we have

$$\begin{aligned} P\left(\max_{k,q,l} \left|[\hat{\mathbf{R}}_{zzk} - \mathbf{R}_{0k}]_{ql}\right| > \frac{\delta}{4}\right) &\leq P_{tb} \\ &= 4M_n(2p_n)^2 \exp\left(-\frac{K_n \delta^2 / 16}{3200 \max_{k,i}([\mathbf{R}_{0k}]_{ii}^2)}\right) \end{aligned} \quad (28)$$

$$= 16M_n p_n^2 \exp\left(-\frac{K_n \delta^2 / 16}{800 \max_{k,i}([\mathbf{S}_{0k}]_{ii}^2)}\right), \quad (29)$$

for all $\delta/4 \in (0, 40 \max_{k,i}([\mathbf{R}_{0k}]_{ii}))$, equivalently, for all $\delta \in (0, c_*^{-1})$, where in the last step above we have used (27), and

$$c_* = (80 \max_{k,i}([\mathbf{S}_{0k}]_{ii}))^{-1}. \quad (30)$$

It follows from (25)-(26) that

$$\begin{aligned} \max_{k,q,l} \left|[\hat{\mathbf{S}}_k - \mathbf{S}_{0k}]_{ql}\right| &\leq \max_{k,q,l} \left|[\hat{\mathbf{R}}_{d_{kr} d_{kr}} - \mathbf{R}_{d_{kr} d_{kr}}]_{ql}\right| \\ &+ \max_{k,q,l} \left|[\hat{\mathbf{R}}_{d_{kr} d_{ki}} - \mathbf{R}_{d_{kr} d_{ki}}]_{ql}\right| + \max_{k,q,l} \left|[\hat{\mathbf{R}}_{d_{ki} d_{kr}} - \mathbf{R}_{d_{ki} d_{kr}}]_{ql}\right| \\ &+ \max_{k,q,l} \left|[\hat{\mathbf{R}}_{d_{ki} d_{ki}} - \mathbf{R}_{d_{ki} d_{ki}}]_{ql}\right| \leq 4 \max_{k,q,l} \left|[\hat{\mathbf{R}}_{zzk} - \mathbf{R}_{0k}]_{ql}\right| \end{aligned}$$

implying

$$\left\{ \max_{k,q,l} \left|[\hat{\mathbf{R}}_{zzk} - \mathbf{R}_{0k}]_{ql}\right| \leq \frac{\delta}{4} \right\} \subseteq \left\{ \max_{k,q,l} \left|[\hat{\mathbf{S}}_k - \mathbf{S}_{0k}]_{ql}\right| \leq \delta \right\}. \quad (31)$$

Hence, using (29), we have

$$\begin{aligned} P\left(\max_{k,q,l} \left|[\hat{\mathbf{S}}_k - \mathbf{S}_{0k}]_{ql}\right| \leq \delta\right) &\geq P\left(\max_{k,q,l} \left|[\hat{\mathbf{R}}_{zzk} - \mathbf{R}_{0k}]_{ql}\right| \leq \frac{\delta}{4}\right) \\ &= 1 - P\left(\max_{k,q,l} \left|[\hat{\mathbf{R}}_{zzk} - \mathbf{R}_{0k}]_{ql}\right| > \frac{\delta}{4}\right) \geq 1 - P_{tb}. \end{aligned}$$

Thus

$$P\left(\max_{k,q,l} \left|[\hat{\mathbf{S}}_k - \mathbf{S}_{0k}]_{ql}\right| > \delta\right) \leq P_{tb}. \quad (32)$$

Now with $C'_0 = c_*^{-1} \sqrt{\frac{2 \ln(16p_n^\tau M_n)}{\ln(p_n)}}$, pick δ to satisfy

$$\delta = C'_0 \sqrt{\frac{\ln(p_n)}{K_n}} \in (0, c_*^{-1}) \Rightarrow c_* \delta = c_* C'_0 \sqrt{\frac{\ln(p_n)}{K_n}} \in (0, 1). \quad (33)$$

Let us pick

$$c_* \delta = \sqrt{\frac{2 \ln(16p_n^\tau M_n)}{K_n}} = \sqrt{N_1/K_n} < 1 \quad (34)$$

for $K_n > N_1$, thereby satisfying $\delta \in (0, c_*^{-1})$. Using (30) in (29), we have

$$\begin{aligned} P_{tb} &= 16M_n p_n^2 / \exp(c_*^2 \delta^2 K_n / 2) = 16M_n p_n^2 / \exp(\ln(16p_n^\tau M_n)) \\ &= \frac{1}{p_n^{\tau-2}}. \end{aligned} \quad (35)$$

Finally, since $\max_{\ell,f} [\mathbf{S}_0(f)]_{\ell\ell} \geq \max_{k,i} [\mathbf{S}_{0k}]_{ii}$, hence, $C_0 \geq C'_0$, bound (24) holds. This completes the proof. ■

Lemma 3 deals with a Taylor series expansion using Wirtinger calculus; its proof is omitted for lack of space.

Lemma 3. For $\Phi_k = \Phi_k^H \succ \mathbf{0}$, define a real scalar function

$$c(\Phi_k, \Phi_k^*) = \ln |\Phi_k| + \ln |\Phi_k^*|. \quad (36)$$

Let $\Phi_k = \Phi_{0k} + \Gamma_k$ with $\Phi_{0k} = \Phi_{0k}^H \succ \mathbf{0}$ and $\Gamma_k = \Gamma_k^H$. Then using Wirtinger calculus, the Taylor series expansion of $c(\Phi_k, \Phi_k^*)$ is given by

$$\begin{aligned} c(\Phi_k, \Phi_k^*) &= c(\Phi_{0k}, \Phi_{0k}^*) + \text{tr}(\Phi_{0k}^{-1} \Gamma_k + \Phi_{0k}^{-*} \Gamma_k^*) \\ &\quad - \frac{1}{2} (\text{vec}(\Gamma_k))^H (\Phi_{0k}^{-*} \otimes \Phi_{0k}^{-1}) \text{vec}(\Gamma_k) \\ &\quad - \frac{1}{2} (\text{vec}(\Gamma_k^*))^H (\Phi_{0k}^{-1} \otimes \Phi_{0k}^{-*}) \text{vec}(\Gamma_k^*) + \text{h.o.t.} \end{aligned} \quad (37)$$

where h.o.t. stands for higher-order terms in Γ_k and Γ_k^* . •

Lemma 3 regarding Taylor series expansion immediately leads to Lemma 4 regarding Taylor series with integral remainder, needed to follow the proof of [17] pertaining to the real-valued case.

Lemma 4. With $c(\Phi_k, \Phi_k^*)$ and $\Phi_k = \Phi_{0k} + \Gamma_k$ as in Lemma 3, the Taylor series expansion of $c(\Phi_k, \Phi_k^*)$ in integral remainder form is given by (v is real)

$$\begin{aligned} c(\Phi_k, \Phi_k^*) &= c(\Phi_{0k}, \Phi_{0k}^*) + \text{tr}(\Phi_{0k}^{-1} \Gamma_k + \Phi_{0k}^{-*} \Gamma_k^*) \\ &\quad - g^H(\Gamma_k) \left(\int_0^1 (1-v) \mathbf{H}(\Phi_{0k}, \Gamma_k, v) dv \right) g(\Gamma_k) \end{aligned} \quad (38)$$

where

$$g(\Gamma_k) = \begin{bmatrix} \text{vec}(\Gamma_k) \\ \text{vec}(\Gamma_k^*) \end{bmatrix}, \quad \mathbf{H}(\Phi_{0k}, \Gamma_k, v) = \begin{bmatrix} \mathbf{H}_{11} & \mathbf{0} \\ \mathbf{0} & \mathbf{H}_{22} \end{bmatrix} \quad (39)$$

$$\mathbf{H}_{11} = (\Phi_{0k} + v\Gamma_k)^{-*} \otimes (\Phi_{0k} + v\Gamma_k)^{-1} \quad (40)$$

and

$$\mathbf{H}_{22} = (\Phi_{0k} + v\Gamma_k)^{-1} \otimes (\Phi_{0k} + v\Gamma_k)^{-*} \quad (41)$$

We now turn to the proof of Theorem 1.

Proof of Theorem 1. Let $\Omega = \Omega_0 + \Delta$ where

$$\Delta = [\Gamma_1 \ \Gamma_2 \ \cdots \ \Gamma_{M_n}] \quad (42)$$

$$\Gamma_k = \Phi_k - \Phi_{0k}, \quad k = 1, 2, \dots, M_n, \quad (43)$$

and Φ_k, Φ_{0k} are both Hermitian positive-definite, implying $\Gamma_k = \Gamma_k^H$. Let

$$Q(\Omega) := \mathcal{L}_{SGL}(\Omega) - \mathcal{L}_{SGL}(\Omega_0). \quad (44)$$

The estimate $\hat{\Omega}_\lambda$, denoted by $\hat{\Omega}$ hereafter suppressing dependence upon λ , minimizes $Q(\Omega)$, or equivalently, $\hat{\Delta} = \hat{\Omega} - \Omega_0$ minimizes $G(\Delta) := Q(\Omega_0 + \Delta)$. We will follow the method of proof of [17, Theorem 1] pertaining to real-valued i.i.d. time series. Consider the set

$$\Theta_n(R) := \left\{ \Delta : \Gamma_k = \Gamma_k^H \ \forall k, \|\Delta\|_F = Rr_n \right\} \quad (45)$$

where R and r_n are as in (17) and (18), respectively. Observe that $G(\Delta)$ is a convex function of Δ , and

$$G(\hat{\Delta}) \leq G(\mathbf{0}) = 0. \quad (46)$$

Therefore, if we can show that

$$\inf_{\Delta} \{G(\Delta) : \Delta \in \Theta_n(R)\} > 0, \quad (47)$$

the minimizer $\hat{\Delta}$ must be inside the sphere defined by $\Theta_n(R)$, and hence

$$\|\hat{\Delta}\|_F \leq Rr_n. \quad (48)$$

Using Lemma 4 we rewrite $G(\Delta)$ as

$$G(\Delta) = \sum_{k=1}^{M_n} (A_{1k} + A_{2k} + A_{3k}) + A_4, \quad (49)$$

where, noting that $\Phi_k^{-1} = S_k$,

$$A_{1k} = g^H(\Gamma_k) \left(\int_0^1 (1-v) \mathbf{H}(\Phi_{0k}, \Gamma_k, v) dv \right) g(\Gamma_k), \quad (50)$$

$$A_{2k} = \text{tr} \left((\hat{S} - S_{0k}) \Gamma_k + (\hat{S} - S_{0k})^* \Gamma_k^* \right), \quad (51)$$

$$A_{3k} = \bar{\lambda}_1 (\|\Phi_{0k}^- + \Gamma_k^- \|_1 - \|\Phi_{0k}^- \|_1), \quad (52)$$

$$A_4 = \bar{\lambda}_2 \sum_{i \neq j}^{p_n} (\|\Omega_0^{(ij)} + \Delta^{(ij)}\|_F - \|\Omega_0^{(ij)}\|_F), \quad (53)$$

$$\Omega_0^{(ij)} := [[\Phi_{01}]_{ij} \ \cdots \ [\Phi_{0M_n}]_{ij}]^T \in \mathbb{C}^{M_n}, \quad (54)$$

$$\Delta^{(ij)} := [[\Gamma_1]_{ij} \ \cdots \ [\Gamma_{M_n}]_{ij}]^T \in \mathbb{C}^{M_n}. \quad (55)$$

Define

$$d_{1n} := \sqrt{\frac{\ln(p_n)}{K_n}}, \quad d_{2n} := d_{1n} \sqrt{p_n + s_{n0}}. \quad (56)$$

Similar to proof of [18, Theorem 1] (see (42) therein), we deduce that

$$A_{1k} \geq \|\Gamma_k\|_F^2 (\beta_{\min}^{-1} + \|\Gamma_k\|_F)^{-2}. \quad (57)$$

But $\|\Gamma_k\|_F \leq \|\Delta\|_F = Rr_n$. Hence

$$A_{1k} \geq \|\Gamma_k\|_F^2 (\beta_{\min}^{-1} + Rr_n)^{-2} \quad (58)$$

and

$$A_1 = \sum_{k=1}^{M_n} A_{1k} \geq \frac{\sum_{k=1}^{M_n} \|\Gamma_k\|_F^2}{(\beta_{\min}^{-1} + Rr_n)^2} = \frac{\|\Delta\|_F^2}{(\beta_{\min}^{-1} + Rr_n)^2} \quad (59)$$

Similar to the proof of [18, Theorem 1] (see (48) therein), we deduce that w.h.p.

$$|A_{2k}| \leq 2C_0 (\|\Gamma_k^- \|_1 d_{1n} + \|\Gamma_k^+ \|_F d_{2n}). \quad (60)$$

Hence with $A_2 = \sum_{k=1}^{M_n} A_{2k}$,

$$|A_2| \leq \sum_{k=1}^{M_n} |A_{2k}| \leq 2C_0 \sum_{k=1}^{M_n} (d_{1n} \|\Gamma_k^- \|_1 + d_{2n} \|\Gamma_k^+ \|_F). \quad (61)$$

We now derive an alternative bound on A_2 . We have w.h.p.

$$|A_2| \leq 2 \sum_{i,j=1}^{p_n} \sum_{k=1}^{M_n} |[(\hat{S} - S_{0k})_{ij}] \cdot [\Gamma_k]_{ij}| \quad (62)$$

$$\leq 2C_0 d_{1n} \sum_{i,j=1}^{p_n} \sum_{k=1}^{M_n} |[\Gamma_k]_{ij}| \quad (63)$$

$$\leq 2C_0 d_{1n} \sum_{i,j=1}^{p_n} (\sqrt{M_n} \|\Delta^{(ij)}\|_F) \quad (64)$$

$$= 2\sqrt{M_n} C_0 d_{1n} (\|\tilde{\Delta}^- \|_1 + \|\tilde{\Delta}^+ \|_1) \quad (65)$$

where $\tilde{\Delta} \in \mathbb{R}^{p_n \times p_n}$ has its (i, j) th element $\tilde{\Delta}_{ij} = \|\Delta^{(ij)}\|_F$.

We now bound A_{3k} . Let \mathcal{E}_0^c denote the complement of \mathcal{E}_0 , given by $\mathcal{E}_0^c = \{(i, j) : [S_0^{-1}(f)]_{ij} \equiv 0, i \neq j, 0 \leq f \leq 0.5\}$. For an index set \mathcal{B} and a matrix $\mathbf{C} \in \mathbb{C}^{p \times p}$, we write $\mathbf{C}_{\mathcal{B}}$ to denote a matrix in $\mathbb{C}^{p \times p}$ such that $[\mathbf{C}_{\mathcal{B}}]_{ij} = C_{ij}$ if $(i, j) \in \mathcal{B}$, and $[\mathbf{C}_{\mathcal{B}}]_{ij} = 0$ if $(i, j) \notin \mathcal{B}$. Then $\Gamma_k^- = \Gamma_{k\mathcal{E}_0}^- + \Gamma_{k\mathcal{E}_0^c}^-$, and $\|\Gamma_k^- \|_1 = \|\Gamma_{k\mathcal{E}_0}^- \|_1 + \|\Gamma_{k\mathcal{E}_0^c}^- \|_1$. We have

$$\begin{aligned} A_{3k} &= \bar{\lambda}_1 (\|\Phi_{0k}^- + \Gamma_k^- \|_1 - \|\Phi_{0k}^- \|_1) \\ &= \bar{\lambda}_1 (\|\Phi_{0k}^- + \Gamma_{k\mathcal{E}_0}^- \|_1 + \|\Gamma_{k\mathcal{E}_0^c}^- \|_1 - \|\Phi_{0k}^- \|_1) \\ &\geq \bar{\lambda}_1 (\|\Gamma_{k\mathcal{E}_0}^- \|_1 - \|\Gamma_{k\mathcal{E}_0}^- \|_1) \end{aligned} \quad (66)$$

leading to $(A_3 = \sum_{k=1}^{M_n} A_{3k})$

$$A_3 \geq \bar{\lambda}_1 \sum_{k=1}^{M_n} (\|\Gamma_{k\varepsilon_0}^-\|_1 - \|\Gamma_{k\varepsilon_0}^-\|_1). \quad (67)$$

Similarly,

$$A_4 \geq \bar{\lambda}_2 (\|\tilde{\Delta}_{\varepsilon_0}^-\|_1 - \|\tilde{\Delta}_{\varepsilon_0}^-\|_1). \quad (68)$$

By Cauchy-Schwartz inequality, $\|\Gamma_{k\varepsilon_0}^-\|_1 \leq \sqrt{s_{n0}} \|\Gamma_{k\varepsilon_0}^-\|_F \leq \sqrt{s_{n0}} \|\Delta\|_F$, hence

$$\sum_{k=1}^{M_n} \|\Gamma_{k\varepsilon_0}^-\|_1 \leq \sqrt{M_n s_{n0}} \|\Delta\|_F. \quad (69)$$

Set $\|\Gamma_k^-\|_1 = \|\Gamma_{k\varepsilon_0}^-\|_1 + \|\Gamma_{k\varepsilon_0^c}^-\|_1$ in A_2 of (61) to deduce that w.h.p.

$$\begin{aligned} \alpha A_2 + A_3 &\geq \alpha(\lambda_n - 2C_0 d_{1n}) \sum_{k=1}^{M_n} \|\Gamma_{k\varepsilon_0^c}^-\|_1 \\ &\quad - \alpha(2C_0 d_{1n} + \lambda_n) \sum_{k=1}^{M_n} \|\Gamma_{k\varepsilon_0}^-\|_1 - \alpha 2C_0 d_{2n} \sum_{k=1}^{M_n} \|\Gamma_k^+\|_F \\ &\geq -\alpha \left((2C_0 d_{1n} + \lambda_n) \sqrt{s_{n0}} - 2C_0 d_{2n} \right) \sqrt{M_n} \|\Delta\|_F \end{aligned} \quad (70)$$

where we have used the fact that $\lambda_n \geq 2C_0 d_{1n}$ and $\sum_{k=1}^{M_n} \|\Gamma_k^+\|_F \leq \sqrt{M_n} \|\Delta\|_F$. Now use A_2 of (65) to deduce that w.h.p.

$$\begin{aligned} (1-\alpha)A_2 + A_4 &\geq (1-\alpha) \left((\lambda_n - 2C_0 \sqrt{M_n} d_{1n}) \|\tilde{\Delta}_{\varepsilon_0}^-\|_1 \right. \\ &\quad \left. - (2C_0 \sqrt{M_n} d_{1n} + \lambda_n) \|\tilde{\Delta}_{\varepsilon_0}^-\|_1 - 2C_0 \sqrt{M_n} p_n d_{1n} \|\Delta\|_F \right) \\ &\geq -(1-\alpha) \|\Delta\|_F \left(\lambda_n \sqrt{s_{n0}} + 2C_0 \sqrt{M_n} d_{1n} (\sqrt{s_{n0}} + \sqrt{p_n}) \right) \end{aligned} \quad (71)$$

where we have used the facts that $\lambda_n \geq 2C_0 \sqrt{M_n} d_{1n}$, and $\|\tilde{\Delta}_{\varepsilon_0}^-\|_1 \leq \sqrt{s_{n0}} \|\tilde{\Delta}_{\varepsilon_0}^-\|_F \leq \sqrt{s_{n0}} \|\Delta\|_F$ (by Cauchy-Schwartz inequality).

Since $r_n = \sqrt{M_n} d_{2n} > \sqrt{M_n s_{n0}} d_{1n}$, w.h.p. we have

$$\begin{aligned} A_2 + A_3 + A_4 &\geq -\|\Delta\|_F \left(\alpha(4C_0 r_n + \lambda_n \sqrt{M_n s_{n0}}) \right. \\ &\quad \left. + (1-\alpha)(\lambda_n \sqrt{s_{n0}} + 4C_0 r_n) \right) \\ &\geq -\|\Delta\|_F \left(4C_0 r_n + \lambda_n \sqrt{s_{n0}} (\alpha \sqrt{M_n} + (1-\alpha)) \right) \\ &\geq -\|\Delta\|_F \left((4+C_1)C_0 r_n \right) \end{aligned} \quad (72)$$

where we have used the fact that, by (18) and (21), $\lambda_n \sqrt{s_{n0}} (\alpha \sqrt{M_n} + (1-\alpha)) \leq C_1 C_0 r_n$. Using (49), (59) and (72), and $\|\Delta\|_F = R r_n$, we have w.h.p.

$$G(\Delta) \geq \|\Delta\|_F^2 \left[(\beta_{\min}^{-1} + R r_n)^{-2} - \frac{(4+C_1)C_0}{R} \right]. \quad (73)$$

For $n \geq N_2$, if we pick R as specified in (17), we obtain $R r_n \leq R r_{N_2} \leq \delta_1 / \beta_{\min}$. Then

$$\begin{aligned} \frac{1}{(\beta_{\min}^{-1} + R r_n)^2} &\geq \frac{\beta_{\min}^2}{(1+\delta_1)^2} = \frac{(4+C_1+\delta_2)C_0}{R} \\ &> \frac{(4+C_1)C_0}{R}, \end{aligned}$$

implying $G(\Delta) > 0$ w.h.p. This proves the desired result. \blacksquare

V. CONCLUSIONS

Graphical modeling of dependent Gaussian time series was considered. A sparse-group lasso-based frequency-domain formulation of the problem has previously been investigated in [11] where the objective was to estimate the inverse power spectral density (PSD) of the data via optimization of a sparse-group lasso based penalized log-likelihood cost function that is formulated in the frequency-domain. The graphical model is then inferred from the estimated inverse PSD. In this paper we established sufficient conditions for consistency of the inverse PSD estimator resulting from the sparse-group graphical lasso-based approach.

REFERENCES

- [1] S.L. Lauritzen, *Graphical models*. Oxford, UK: Oxford Univ. Press, 1996.
- [2] 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.
- [3] R. Dahlhaus, "Graphical interaction models for multivariate time series," *Metrika*, vol. 51, pp. 157-172, 2000.
- [4] M. Eichler, "Graphical modelling of multivariate time series," *Probability Theory and Related Fields*, vol. 153, issue 1-2, pp. 233-268, June 2012.
- [5] 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.
- [6] N. Friedman, "Inferring cellular networks using probabilistic graphical models," *Science*, vol 303, pp. 799-805, 2004.
- [7] S.L. Lauritzen and N.A. Sheehan, "Graphical models for genetic analyses," *Statistical Science*, vol. 18, pp. 489-514, 2003.
- [8] N. Meinshausen and P. Bühlmann, "High-dimensional graphs and variable selection with the Lasso," *Ann. Statist.*, vol. 34, no. 3, pp. 1436-1462, 2006.
- [9] K. Mohan, M.J.Y. Chung, S. Han, D. Witten, S.I. Lee and M. Fazel, "Structured learning of Gaussian graphical models," in *Proc. NIPS 2012*, Lake Tahoe, Dec. 2012, pp. 620-628.
- [10] 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.
- [11] J.K. Tugnait, "Graphical modeling of high-dimensional time series," in *Proc. 52nd Asilomar Conference on Signals, Systems and Computers*, pp. 840-844, Pacific Grove, CA, Oct. 29 - Oct. 31, 2018.
- [12] 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.
- [13] 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.
- [14] D.R. Brillinger, *Time Series: Data Analysis and Theory*, Expanded edition. New York: McGraw Hill, 1981.
- [15] P.J. Schreier and L.L. Scharf, *Statistical Signal Processing of Complex-Valued Data*, Cambridge, UK: Cambridge Univ. Press, 2010.
- [16] N. Simon, J. Friedman, T. Hastie and R. Tibshirani, "A sparse-group lasso," *J. Computational Graphical Statistics*, vol. 22, pp. 231-245, 2013.
- [17] 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.
- [18] J.K. Tugnait, "On sparse complex Gaussian graphical model selection," in *Proc. 2019 IEEE Intern. Workshop on Machine Learning for Signal Processing (MLSP 2019)*, Pittsburgh, PA, Oct. 13-16, 2019.
- [19] 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.