

Calibrated Precision Matrix Estimation for High-Dimensional Elliptical Distributions

Tuo Zhao and Han Liu

Abstract—We propose a semiparametric method for estimating a precision matrix of high-dimensional elliptical distributions. Unlike most existing methods, our method naturally handles heavy tailness and conducts parameter estimation under a calibration framework, thus achieves improved theoretical rates of convergence and finite sample performance on heavy-tail applications. We further demonstrate the performance of the proposed method using thorough numerical experiments.

Index Terms—Precision matrix, calibrated estimation, elliptical distribution, heavy-tailness, semiparametric model.

I. INTRODUCTION

WE CONSIDER the problem of precision matrix estimation. Let $X = (X_1, \dots, X_d)^T$ be a d -dimensional random vector with mean $\mu \in \mathbb{R}^d$ and covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$, where $\Sigma_{kj} = \mathbb{E}X_k X_j - \mathbb{E}X_k \mathbb{E}X_j$. We want to estimate the precision matrix $\Omega = \Sigma^{-1}$ based on n independent observations. In this paper we focus on high dimensional settings where $d/n \rightarrow \infty$. To handle the curse of dimensionality, we assume that Ω is sparse (i.e., many off-diagonal entries of Ω are zero).

A popular statistical model for precision matrix estimation is multivariate Gaussian, i.e., $X \sim N(\mu, \Sigma)$. Under Gaussian models, sparse precision matrix encodes the conditional independence relationship of the random variables [8], [21], which has motivated numerous applications in different research areas [3], [15], [36]. In the past decade, many precision matrix estimation methods have been proposed for Gaussian distributions. For more details, let $x_1, \dots, x_n \in \mathbb{R}^d$ be n independent observations of X , we define the sample covariance matrix as

$$S = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T, \quad (1)$$

Manuscript received July 29, 2013; revised May 9, 2014; accepted September 11, 2014. Date of publication September 30, 2014; date of current version November 18, 2014. This work was supported in part by the National Science Foundation under Grant IIS1408910 and Grant IIS1332109 and in part by the National Institutes of Health under Grant R01MH102339, Grant R01GM083084, and Grant R01HG06841. This paper was presented at the 27th Annual Conference on Neural Information Processing Systems in 2013.

T. Zhao is with the Department of Operations Research and Financial Engineering, Princeton University, Princeton, NJ 08544 USA, and also with the Department of Computer Science, Johns Hopkins University, Baltimore, MD 21218 USA (e-mail: tour@cs.jhu.edu).

H. Liu is with the Department of Operations Research and Financial Engineering, Princeton University, Princeton, NJ 08544 USA (e-mail: hanliu@princeton.edu).

Communicated by N. Cesa-Bianchi, Associate Editor for Pattern Recognition, Statistical Learning and Inference.

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TIT.2014.2360980

where $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$. [1], [11], [38] propose the penalized Gaussian log-likelihood method named graphical lasso (GLASSO), which solves

$$\hat{\Omega} = \underset{\Omega}{\operatorname{argmin}} -\log |\Omega| + \operatorname{tr}(S\Omega) + \lambda \sum_{k,j} |\Omega_{kj}|, \quad (2)$$

where $\lambda > 0$ is a regularization parameter for controlling the bias-variance tradeoff. In another line of research, [5], [37] propose pseudo-likelihood methods to estimate the precision matrix. Their methods adopt a column-by-column estimation scheme and are more amenable to theoretical analysis. More specifically, given a matrix $A \in \mathbb{R}^{d \times d}$, let $A_{*j} = (A_{1j}, \dots, A_{dj})^T$ denote the j th column of A , we define $\|A_{*j}\|_1 = \sum_k |A_{kj}|$ and $\|A_{*j}\|_\infty = \max_k |A_{kj}|$. [5] propose the CLIME estimator, which solves

$$\begin{aligned} \hat{\Omega}_{*j} &= \underset{\Omega_{*j}}{\operatorname{argmin}} \|\Omega_{*j}\|_1 \\ \text{s.t. } &\|\Sigma \Omega_{*j} - \mathbf{I}_{*j}\|_\infty \leq \lambda, \quad \forall j = 1, \dots, d, \end{aligned} \quad (3)$$

to estimate the j th column of the precision matrix. Moreover, let $\|A\|_1 = \max_j \|A_{*j}\|_1$ be the matrix ℓ_1 norm of A , and $\|A\|_2$ be the largest singular value of A , (i.e., the spectral norm of A), [5] show that if we choose

$$\lambda \asymp \|\Omega\|_1 \cdot \sqrt{\frac{\log d}{n}}, \quad (4)$$

the CLIME estimator in (3) attains the rates of convergence

$$\|\hat{\Omega} - \Omega\|_p = O_P \left(\|\Omega\|_1^2 \cdot s \sqrt{\frac{\log d}{n}} \right), \quad (5)$$

where $s = \max_j \sum_k \mathbf{1}(\Omega_{kj} \neq 0)$, and $p = 1, 2$. Scalable software packages for GLASSO and CLIME have been developed which scale to thousands of dimensions [16], [22], [40].

Though significant progress has been for estimating Gaussian graphical models, most existing methods have two drawbacks: (i) They generally require the underlying distribution to be light-tailed [5], [7]. When this assumption is violated, these sample covariance matrix-based methods may have poor performance. (ii) They generally use the same tuning parameter to regularize the estimation, which is not adaptive to the individual sparseness of each column (More details will be provided in §III.B) and may lead to inferior finite sample performance. In another word, the regularization for estimating different columns of the precision matrix is not calibrated.

To overcome the above drawbacks, we propose a new sparse precision matrix estimation method, named EPIC (Estimating Precision matrix with Calibration), which simultaneously

handles data heavy-tailness and conducts calibrated estimation. To relax the tail conditions, we adopt a combination of the rank-based transformed Kendall's tau estimator and Catoni's M-estimator [7], [18]. Such a semiparametric combination has shown better statistical properties than those of the sample covariance matrix for the heavy-tailed elliptical distributions [6], [7], [10], [17]. We will explain more details in § II and § IV. To calibrate the parameter estimation, we exploit a new framework proposed by [12]. Under this framework, the optimal tuning parameter does not depend on any unknown quantity of the data distribution, thus the EPIC estimator is tuning insensitive [25]. Computationally, the EPIC estimator is formulated as a convex program, which can be efficiently solved by the parametric simplex method [34]. Theoretically, we show that the EPIC estimator attains improved rates of convergence than the one in (5) under mild conditions. Numerical experiments on both simulated and real datasets show that the EPIC method outperforms existing precision matrix estimation methods.

The rest of this paper is organized as follows: In §II, we briefly review the elliptical family; In §III, we describe the proposed method and derive the computational algorithm; In §IV, we analyze the statistical properties of the EPIC estimator; In §V and §VI, we conduct numerical experiments on both simulated and real datasets to illustrate the effectiveness of the proposed method; In §VII, we discuss other related precision matrix estimation methods and compare them with our method [23]–[25].

II. BACKGROUND

We start with some notations. Let $\mathbf{v} = (v_1, \dots, v_d)^T \in \mathbb{R}^d$ be a vector, we define vector norms: $\|\mathbf{v}\|_1 = \sum_{j=1}^d |v_j|$, $\|\mathbf{v}\|_2^2 = \sum_{j=1}^d v_j^2$, $\|\mathbf{v}\|_\infty = \max_{1 \leq j \leq d} |v_j|$. Let \mathcal{S} be a subspace of \mathbb{R}^d , we use $\mathbf{v}_\mathcal{S}$ to denote the projection of \mathbf{v} onto \mathcal{S} : $\mathbf{v}_\mathcal{S} = \operatorname{argmin}_{\mathbf{u} \in \mathcal{S}} \|\mathbf{u} - \mathbf{v}\|_2^2$. We also define the orthogonal complement of \mathcal{S} as $\mathcal{S}^\perp = \{\mathbf{u} \in \mathbb{R}^d | \mathbf{u}^T \mathbf{v} = 0, \text{ for any } \mathbf{v} \in \mathcal{S}\}$. Given a matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$, let $\mathbf{A}_{*j} = (\mathbf{A}_{1j}, \dots, \mathbf{A}_{dj})^T$ and $\mathbf{A}_{k*} = (\mathbf{A}_{k1}, \dots, \mathbf{A}_{kd})^T$ denote the j^{th} column and k^{th} row of \mathbf{A} in vector forms, we define matrix norms: $\|\mathbf{A}\|_1 = \max_j \|\mathbf{A}_{*j}\|_1$, $\|\mathbf{A}\|_2 = \psi_{\max}(\mathbf{A})$, $\|\mathbf{A}\|_\infty = \max_k \|\mathbf{A}_{k*}\|_1$, $\|\mathbf{A}\|_F^2 = \sum_j \|\mathbf{A}_{*j}\|_2^2$, $\|\mathbf{A}\|_{\max} = \max_j \|\mathbf{A}_{*j}\|_\infty$, where $\psi_{\max}(\mathbf{A})$ is the largest singular value of \mathbf{A} . We use $\Lambda_{\max}(\mathbf{A})$ and $\Lambda_{\min}(\mathbf{A})$ to denote the largest and smallest eigenvalues of \mathbf{A} . Moreover, we define the projection of \mathbf{A}_{*j} onto \mathcal{S} as $\mathbf{A}_{\mathcal{S}j} = \operatorname{argmin}_{\mathbf{u} \in \mathcal{S}} \|\mathbf{u} - \mathbf{A}_{*j}\|_2^2$.

We then briefly review the elliptical family, which has the following definition.

Definition 2.1 ([10]): Given $\boldsymbol{\mu} \in \mathbb{R}^d$ and a symmetric positive semidefinite matrix $\boldsymbol{\Sigma}$ with $\operatorname{rank}(\boldsymbol{\Sigma}) = r \leq d$, we say that a d -dimensional random vector $\mathbf{X} = (X_1, \dots, X_d)^T$ follows an elliptical distribution with parameter $\boldsymbol{\mu}$, $\boldsymbol{\Sigma}$, and $\boldsymbol{\Sigma}$ denoted by

$$\mathbf{X} \sim EC(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}), \quad (6)$$

if \mathbf{X} has a stochastic representation

$$\mathbf{X} \stackrel{d}{=} \boldsymbol{\mu} + \boldsymbol{\Sigma} \mathbf{A} \mathbf{U} \quad (7)$$

where $\xi \geq 0$ is a continuous random variable independent of \mathbf{U} . Here $\mathbf{U} \in \mathbb{S}^{r-1}$ is uniformly distributed on the unit sphere in \mathbb{R}^r , and $\boldsymbol{\Sigma} = \mathbf{A} \mathbf{A}^T$.

Note that \mathbf{A} and ξ in (7) can be properly rescaled without changing the distribution. Thus existing literature usually imposes an additional constraint $\|\boldsymbol{\Sigma}\|_{\max} = 1$ to make the distribution identifiable [10]. However, such a constraint does not necessarily make $\boldsymbol{\Sigma}$ the covariance matrix of \mathbf{X} . Since we are interested in estimating the precision matrix in this paper, we require $\mathbb{E}(\xi^2) < \infty$ and $\operatorname{rank}(\boldsymbol{\Sigma}) = d$ such that the precision matrix of the elliptical distribution exists. Under this assumption, we use an alternative constraint $\mathbb{E}(\xi^2) = d$, which not only makes the distribution identifiable but also has $\boldsymbol{\Sigma}$ defined as the conventional covariance matrix (e.g., as in the Gaussian distribution).

Remark 1: $\boldsymbol{\Sigma}$ can be factorized as $\boldsymbol{\Sigma} = \boldsymbol{\Theta} \mathbf{Z} \boldsymbol{\Theta}$, where \mathbf{Z} is the Pearson correlation matrix, and $\boldsymbol{\Theta} = \operatorname{diag}(\theta_1, \dots, \theta_d)$ with θ_j as the standard deviation of X_j . Since $\boldsymbol{\Theta}$ is a diagonal matrix, we can rewrite the precision matrix $\boldsymbol{\Omega}$ as $\boldsymbol{\Omega} = \boldsymbol{\Theta}^{-1} \boldsymbol{\Gamma} \boldsymbol{\Theta}^{-1}$, where $\boldsymbol{\Gamma} = \mathbf{Z}^{-1}$ is the inverse correlation matrix.

Remark 2: As a generalization of the Gaussian family, the elliptical family has been widely applied to many research areas such as dimensionality reduction [19], portfolio theory [14], and data visualization [33]. Many of these applications rely on an effective estimator of the precision matrix for elliptical distributions.

III. METHOD

Motivated by the above discussion, the EPIC method has three steps: We first use the transformed Kendall's tau estimator and Catoni's M-estimator to obtain $\widehat{\mathbf{Z}}$ and $\widehat{\boldsymbol{\Theta}}$ respectively; We then plug $\widehat{\mathbf{Z}}$ into a calibrated inverse correlation matrix estimation procedure to obtain $\widehat{\boldsymbol{\Gamma}}$; At last we assemble $\widehat{\boldsymbol{\Gamma}}$ and $\widehat{\boldsymbol{\Theta}}$ to obtain $\widehat{\boldsymbol{\Omega}}$. We explain more details about these three steps in the following subsections.

A. Correlation Matrix and Standard Deviation Estimation

To estimate \mathbf{Z} , we adopt the transformed Kendall's tau estimator proposed in [10] and [23]. More specifically, we define a population version of the Kendall's tau statistic between X_j and X_k as follows,

$$\begin{aligned} \tau_{kj} &= \mathbb{P}((X_j - \widetilde{X}_j)(X_k - \widetilde{X}_k) > 0) \\ &\quad - \mathbb{P}((X_j - \widetilde{X}_j)(X_k - \widetilde{X}_k) < 0), \end{aligned}$$

where \widetilde{X}_j and \widetilde{X}_k are independent copies of X_j and X_k respectively. For elliptical distributions, [10], [23] show that \mathbf{Z}_{kj} 's and τ_{kj} 's have the following relationship

$$\mathbf{Z} = [\mathbf{Z}_{kj}] = \left[\sin\left(\frac{\pi}{2} \tau_{kj}\right) \right]. \quad (8)$$

Therefore given $\mathbf{x}_1, \dots, \mathbf{x}_n$ be n independent observations of \mathbf{X} , where $\mathbf{x}_i = (x_{i1}, \dots, x_{id})^T$, we first calculate a sample version of the Kendall's tau statistic between X_j and X_k by

$$\widehat{\tau}_{kj} = \frac{2 \sum_{i < i'} \operatorname{sign}((x_{ik} - x_{i'k})(x_{ij} - x_{i'j}))}{n(n-1)}$$

for all $k \neq j$, and 1 otherwise. We then obtain a correlation matrix estimator by the same entrywise transformation as (8),

$$\widehat{\mathbf{Z}} = [\widehat{\mathbf{Z}}_{kj}] = \left[\sin \left(\frac{\pi}{2} \widehat{\tau}_{kj} \right) \right]. \quad (9)$$

To estimate Θ , we exploit the Catoni's M-estimator proposed in [7]. For heavy-tailed distributions, [7] show that the Catoni's M-estimator has better theoretical and empirical performance than the sample moment-based estimator. In particular, let $\psi(t) = \text{sign}(t) \cdot \log(1 + |t| + t^2/2)$ be a univariate function where $\text{sign}(0) = 0$. Let $\widehat{\mu}_j$ and \widehat{m}_j be the estimator of $\mathbb{E}X_j$ and $\mathbb{E}X_j^2$ respectively which solve the following two equations:

$$\sum_{i=1}^n \psi \left((x_{ij} - \mu_j) \sqrt{\frac{2}{nK_{\max}}} \right) = 0, \quad (10)$$

$$\sum_{i=1}^n \psi \left((x_{ij}^2 - m_j) \sqrt{\frac{2}{nK_{\max}}} \right) = 0. \quad (11)$$

Here K_{\max} is a preset upper bound of $\max_j \text{Var}(X_j)$ and $\max_j \text{Var}(X_j^2)$. [7] shows that the solutions to (10) and (11) must exist and can be efficiently solved by the Newton-Raphson algorithm [31]. Once we obtain \widehat{m}_j and $\widehat{\mu}_j$, we estimate the marginal standard deviation θ_j by

$$\widehat{\theta}_j = \sqrt{\max\{\widehat{m}_j - \widehat{\mu}_j^2, K_{\min}\}}, \quad (12)$$

where K_{\min} is a preset lower bound of $\min_j \theta_j^2$.

Remark 3: We choose the combination of the transformed Kendall's tau estimator and Catoni's M-estimator instead of sample covariance matrix, because we are handling heavy-tailed elliptical distributions. For light-tailed distributions (e.g. Gaussian distribution), we can still use the sample correlation matrix and sample standard deviation to estimate the \mathbf{Z} and Θ . The extension of our proposed methodology and theory is straightforward. See more details in §IV.

B. Calibrated Inverse Correlation Matrix Estimation

We then plug the transformed Kendall's tau estimator $\widehat{\mathbf{Z}}$ into the following convex program,

$$\begin{aligned} (\widehat{\Gamma}_{*j}, \widehat{\tau}_j) &= \underset{\Gamma_{*j}, \tau_j}{\text{argmin}} \quad \|\Gamma_{*j}\|_1 + c\tau_j \\ \text{s.t.} \quad &\|\widehat{\mathbf{Z}}\Gamma_{*j} - \mathbf{I}_{*j}\|_{\infty} \leq \lambda\tau_j, \quad \|\Gamma_{*j}\|_1 \leq \tau_j, \end{aligned} \quad (13)$$

for all $j = 1, \dots, d$, where c can be any constant between 0 and 1 (e.g., $c = 0.5$). Here τ_j serves as an auxiliary variable to calibrate the regularization [12], [32]. Both the objective function and constraints in (13) contain τ_j to prevent from choosing τ_j either too large or too small.

To gain more intuition of the formulation of (13), we first consider estimating the j^{th} column of the inverse correlation matrix using the CLIME method in a regularization form as follows,

$$\widehat{\Gamma}_{*j} = \underset{\Gamma_{*j}}{\text{argmin}} \quad \|\Gamma_{*j}\|_1 + \nu \|\widehat{\mathbf{Z}}\Gamma_{*j} - \mathbf{I}_{*j}\|_{\infty}, \quad (14)$$

where $\nu > 0$ is the regularization parameter. The next proposition presents an alternative formulation of (14).

Proposition III.1: The following optimization problem

$$\begin{aligned} (\widehat{\Gamma}_{*j}, \widehat{\tau}_j) &= \underset{\Gamma_{*j}, \tau_j}{\text{argmin}} \quad \|\Gamma_{*j}\|_1 + c\tau_j \\ \text{s.t.} \quad &\|\widehat{\mathbf{Z}}\Gamma_{*j} - \mathbf{I}_{*j}\|_{\infty} \leq \frac{c}{\nu}\tau_j. \end{aligned} \quad (15)$$

has the same solution as (14).

The proof of Proposition III.1 is provided in Appendix A. If we set $\nu/c = \lambda$, then the only difference between (13) and (15) is that (13) contains a constraint $\|\Gamma_{*j}\|_1 \leq \tau_j$. Due to the complementary slackness, this additional constraint encourages the regularization $\lambda\tau_j$ to be proportional to the ℓ_1 norm of the j^{th} column (weak sparseness). From the theoretical analysis in §IV, we see that the regularization is calibrated in this way.

In the rest of this subsection, we omit the index j in (13) for notational simplicity. We denote Γ_{*j} , \mathbf{I}_{*j} , and τ_j by $\boldsymbol{\gamma}$, \mathbf{e} , and τ respectively. By reparametrizing $\boldsymbol{\gamma} = \boldsymbol{\gamma}^+ - \boldsymbol{\gamma}^-$, we can rewrite (13) as the following linear program,

$$\begin{aligned} (\widehat{\boldsymbol{\gamma}}^+, \widehat{\boldsymbol{\gamma}}^-, \widehat{\tau}) &= \underset{\boldsymbol{\gamma}^+, \boldsymbol{\gamma}^-, \tau}{\text{argmin}} \quad \mathbf{1}^T \boldsymbol{\gamma}^+ + \mathbf{1}^T \boldsymbol{\gamma}^- + c\tau \\ \text{s.t.} \quad &\begin{bmatrix} \widehat{\mathbf{Z}} & -\widehat{\mathbf{Z}} & -\lambda \\ -\widehat{\mathbf{Z}} & \widehat{\mathbf{Z}} & -\lambda \\ \mathbf{1}^T & \mathbf{1}^T & -1 \end{bmatrix} \begin{bmatrix} \boldsymbol{\gamma}^+ \\ \boldsymbol{\gamma}^- \\ \tau \end{bmatrix} \leq \begin{bmatrix} \mathbf{e} \\ -\mathbf{e} \\ 0 \end{bmatrix}, \\ &\boldsymbol{\gamma}^+ \geq \mathbf{0}, \quad \boldsymbol{\gamma}^- \geq \mathbf{0}, \quad \tau \geq 0, \end{aligned} \quad (16)$$

where $\lambda = \lambda\mathbf{1}$. Though (16) can be solved by general linear program solvers (e.g. the simplex method as suggested in [5]), these general solvers cannot scale to large problems. In Appendix B, we provide a more efficient parametric simplex method [34], which naturally exploits the underlying sparsity structure, and attains better empirical performance than the simplex method.

C. Symmetric Precision Matrix Estimation

Once we get the inverse correlation matrix estimate $\widehat{\Gamma}$, we estimate the precision matrix by

$$\widetilde{\Omega} = \widehat{\Theta}^{-1} \widehat{\Gamma} \widehat{\Theta}^{-1}.$$

Remark 4: A possible alternative is that we first assemble a covariance matrix estimator

$$\widehat{\mathbf{S}} = \widehat{\Theta} \widehat{\mathbf{Z}} \widehat{\Theta}, \quad (17)$$

then directly estimate Ω by solving

$$\begin{aligned} (\widehat{\Omega}_{*j}, \widehat{\tau}_j) &= \underset{\Omega_{*j}, \tau_j}{\text{argmin}} \quad \|\Omega_{*j}\|_1 + c\tau_j \\ \text{s.t.} \quad &\|\widehat{\mathbf{S}}\Omega_{*j} - \mathbf{I}_{*j}\|_{\infty} \leq \lambda\tau_j, \quad \|\Omega_{*j}\|_1 \leq \tau_j \end{aligned}$$

for all $j = 1, \dots, d$. However, such a direct estimation procedure makes the regularization parameter selection sensitive to marginal variability. See [20], [26], [29] for more discussions of the ensemble rule.

The EPIC method does not guarantee the symmetry of $\widetilde{\Omega}$. To get a symmetric estimate, we take an additional projection procedure to obtain a symmetric estimator

$$\widehat{\Omega} = \underset{\Omega}{\text{argmin}} \quad \|\Omega - \widetilde{\Omega}\|_* \quad \text{s.t.} \quad \Omega = \Omega^T, \quad (18)$$

where $\|\cdot\|_*$ can be the matrix ℓ_1 , Frobenius, or max norm. More details about how to choose a suitable norm will be explained in the next section.

Remark 5: For the Frobenius and max norms, (18) has a closed form solution as follows,

$$\widehat{\mathbf{\Omega}} = \frac{1}{2} (\tilde{\mathbf{\Omega}} + \tilde{\mathbf{\Omega}}^T).$$

For the matrix ℓ_1 norm, see our proposed smoothed proximal gradient algorithm in Appendix C. More details about how to choose a suitable norm will be explained in the next section.

IV. STATISTICAL PROPERTIES

To analyze the statistical properties of the EPIC estimator, we define the following class of sparse symmetric matrices,

$$\mathcal{U}(s, M, \kappa_u) = \left\{ \mathbf{\Gamma} \in \mathbb{R}^{d \times d} \mid \mathbf{\Gamma} \succ 0, \Lambda_{\max}(\mathbf{\Gamma}) \leq \kappa_u, \right. \\ \left. \max_j \sum_k I(\Gamma_{kj} \neq 0) \leq s, \|\mathbf{\Gamma}\|_1 \leq M \right\},$$

where κ_u is a constant, and (s, d, M) may scale with the sample size n . We assume that the following conditions hold:

- (A.1) $\mathbf{\Gamma} \in \mathcal{U}(s, M, \kappa_u)$,
- (A.2) $\theta_{\min} \leq \min_j \theta_j \leq \max_j \theta_j \leq \theta_{\max}$,
- (A.3) $\max_j |\mu_j| \leq \mu_{\max}$, $\max_j \mathbb{E} X_j^4 \leq K$,
- (A.4) $s^2 \log d/n \rightarrow 0$,

where θ_{\max} , θ_{\min} , μ_{\max} , and K are constants.

Remark 6: Condition (A.3) only requires the fourth moment of the distribution to be finite. In contrast, sample covariance-based estimation methods can not achieve such theoretical results. See more details in [5] and [7].

Remark 7: The bounded mean in Condition (A.3) is actually a mild assumption. Existing high dimensional theories (Cai et al. 2011; Yuan, 2010; Rothman et al. 2008) on sparse precision matrix estimation all require the distribution to be light-tailed. For example, there exists some constant K such that $\max_j \mathbb{E}|X_j|^r \leq K < \infty$ for some $r \gg 4$. By Jessen's inequality, we have $(\mathbb{E}|X_j|)^r \leq \mathbb{E}|X_j|^r \leq K < \infty$, which implies that $\max_j \mathbb{E}|X_j| \leq K^{1/r} < \infty$. In another word, they also require $\max_j |\mu_j|$ to be bounded.

Before we proceed with main results, we first present the following important lemma.

Lemma 1: We assume that $\mathbf{X} \sim EC(\boldsymbol{\mu}, \boldsymbol{\xi}, \boldsymbol{\Sigma})$ and (A.2)-(A.4) hold. Let $\widehat{\mathbf{Z}}$ and $\widehat{\theta}_j$ be defined in (9) and (12). There exist universal constants κ_1 and κ_2 such that for large enough n ,

$$\mathbb{P}\left(\max_{k,j} |\widehat{\mathbf{Z}}_{kj} - \mathbf{Z}_{kj}| \leq \kappa_1 \sqrt{\frac{\log d}{n}}\right) \geq 1 - \frac{1}{d}, \quad (19)$$

$$\mathbb{P}\left(\max_j |\widehat{\theta}_j^{-1} - \theta_j^{-1}| \leq \kappa_2 \sqrt{\frac{\log d}{n}}\right) \geq 1 - \frac{2}{d}. \quad (20)$$

The proof of Lemma 1 is provided in Appendix D.

Remark 8: Lemma 1 shows that the transformed Kendall's tau estimator and Catoni's M-estimator possess good concentration properties for heavy-tailed elliptical distributions. That enables us to obtain a consistent precision matrix estimator in high dimensions.

A. Parameter Estimation Consistency

Theorem IV.1 provides the rates of convergence for precision matrix estimation under the matrix ℓ_1 , spectral, and Frobenius norms.

Theorem IV.1: Suppose that $\mathbf{X} \sim EC(\boldsymbol{\mu}, \boldsymbol{\xi}, \boldsymbol{\Sigma})$ and (A.1)-(A.4) hold, if we take $\lambda = \kappa_1 \sqrt{\log d/n}$ and choose the matrix ℓ_1 norm as $\|\cdot\|_*$ in (18), then for large enough n and $p = 1, 2$, there exists a universal constant C_1 such that

$$\mathbb{P}\left(\|\widehat{\mathbf{\Omega}} - \mathbf{\Omega}\|_p \leq C_1 M s \sqrt{\frac{\log d}{n}}\right) \geq 1 - \frac{3}{d}. \quad (21)$$

Moreover, if we choose the Frobenius norm as $\|\cdot\|_*$ in (18), then for large enough n , there exists a universal constant C_2 such that

$$\mathbb{P}\left(\frac{1}{d} \|\widehat{\mathbf{\Omega}} - \mathbf{\Omega}\|_F^2 \leq C_2 M^2 \frac{s \log d}{n}\right) \geq 1 - \frac{3}{d}. \quad (22)$$

The proof of Theorem IV.1 is provided in Appendix E. Note that the rates of convergence obtained in the above theorem are faster than those in [5].

B. Model Selection Consistency

Theorem IV.2 provides the rate of convergence under the elementwise max norm.

Theorem IV.2: Suppose that $\mathbf{X} \sim EC(\boldsymbol{\mu}, \boldsymbol{\xi}, \boldsymbol{\Sigma})$ and (A.1)-(A.4) hold. If we take $\lambda = \kappa_1 \sqrt{\log d/n}$ and choose the max norm for (18), then for large enough n , there exists a universal constant C_3 such that

$$\mathbb{P}\left(\|\widehat{\mathbf{\Omega}} - \mathbf{\Omega}\|_{\max} \leq C_3 M^2 \sqrt{\frac{\log d}{n}}\right) \geq 1 - \frac{3}{d}. \quad (23)$$

Moreover, let $E = \{(k, j) \mid \mathbf{\Omega}_{kj} \neq 0\}$, and $\widehat{E} = \{(k, j) \mid \widehat{\mathbf{\Omega}}_{kj} \neq 0\}$, if there exists large enough constant C_4 such that

$$\min_{(k,j) \in E} |\mathbf{\Omega}_{kj}| \geq C_4 M^2 \sqrt{\frac{\log d}{n}},$$

then we have $\mathbb{P}(E \subseteq \widehat{E}) \rightarrow 1$.

The proof of Theorem IV.2 is provided in Appendix G. The obtained rate of convergence in Theorem IV.2 is comparable to that of [5].

Remark 9: Our selected regularization parameter $\lambda = \kappa_1 \sqrt{\log d/n}$ in Theorems IV.1 and IV.2 does not contain any unknown parameter of the underlying distribution (e.g. $\|\mathbf{\Gamma}\|_1$). Note that κ_1 comes from (19) in Lemma 1. Theoretically we can choose κ_1 as a reasonably large constant without any additional tuning (e.g. $\sqrt{2}\pi$. See more details in [23]). In practice, we found that a fine tuning of κ_1 delivers better finite sample performance.

V. NUMERICAL RESULTS

In this section, we compare the EPIC estimator with several competing estimators including:

- CLIME.RC: We obtain the sparse precision matrix estimator by plugging the covariance matrix estimator $\widehat{\mathbf{S}}$ defined in (17) into (3).

TABLE I

TIMING PERFORMANCE OF DIFFERENT ESTIMATORS ON THE BAND, ERDÖS-RÉNYI, AND SCALE-FREE MODELS (IN SECONDS). THE BASELINE PERFORMANCE IS OBTAINED BY SOLVING THE CLIME.SC METHOD USING THE SIMPLEX METHOD

Model	d	EPIC	GLASSO.RC	CLIME.RC	CLIME.SC	BASELINE
Band	101	0.1561(0.0248)	0.3633(0.0070)	0.1233(0.0057)	0.1701(0.0119)	49.467(1.7862)
	201	1.6622(0.1253)	0.4417(0.0122)	1.5897(0.1249)	1.6085(0.0518)	687.57(23.720)
	401	23.061(0.5777)	1.0864(0.1403)	24.441(1.5344)	25.445(3.8066)	4756.4(170.25)
Erdős-Rényi	101	0.1414(0.0079)	0.3703(0.0072)	0.1309(0.0331)	0.2073(0.0925)	59.775(2.0521)
	201	1.6214(0.5175)	0.4448(0.0164)	1.5992(0.1840)	1.6155(0.2957)	803.51(29.835)
	401	21.722(0.5470)	1.1517(0.0959)	22.795(0.6999)	24.230(3.1871)	4531.7(151.46)
Scale-free	101	0.2245(0.0514)	0.4398(0.0843)	0.1509(0.0054)	0.1871(0.0149)	55.112(1.7109)
	201	1.8682(0.1078)	0.4632(0.0067)	1.5472(0.1350)	1.7235(0.1778)	865.98(31.399)
	401	21.926(0.7112)	1.0093(0.1140)	23.135(1.4318)	25.596(3.3401)	4991.2(202.44)

- **CLIME.SC**: We obtain the sparse precision matrix estimator by plugging the sample covariance matrix estimator \mathbf{S} defined in (1) into (3).
- **GLASSO.RC**: We obtain the sparse precision matrix estimator by plugging the covariance matrix estimator $\hat{\mathbf{S}}$ defined in (17) into (2).

Moreover, (3) is also solved by the parametric simplex method as our proposed EPIC method, and (2) is solved by the block coordinate descent algorithm. All experiments are conducted on a PC with Core i5 3.3GHz CPU and 16GB memory. All programs are coded using C using double precision, and further called from R.

A. Data Generation

We consider three different settings for comparison: (1) $d = 101$; (2) $d = 201$; (3) $d = 401$. We adopt the following three graph generation schemes, as illustrated in Figure 1, to obtain precision matrices:

- **Band**. Each node is assigned an index j with $j = 1, \dots, d$. Two nodes are connected by an edge if the difference between their indices is no larger than 2.
- **Erdős-Rényi**. We set an edge between each pair of nodes with probability $4/d$, independently of the other edges.
- **Scale-free**. The degree distribution of the graph follows a power law. The graph is generated by the preferential attachment mechanism.

The graph begins with an initial chain graph of 10 nodes. New nodes are added to the graph one at a time. Each new node is connected to an existing node with a probability that is proportional to the number of degrees that the existing nodes already have. Formally, the probability p_i that the new node is connected to the i^{th} existing node is $p_i = \frac{k_i}{\sum_j k_j}$ where k_i is the degree of node i .

Let \mathbf{G} be the adjacency matrix of the generated graph, we calculate $\tilde{\mathbf{G}} = [\tilde{\mathbf{G}}_{jk}]$ as

$$\tilde{\mathbf{G}}_{jk} = \tilde{\mathbf{G}}_{kj} = \begin{cases} \mathbf{U}_{kj} & \text{if } \mathbf{G}_{jk} = \mathbf{G}_{kj} = 1 \\ 0 & \text{if } \mathbf{G}_{jk} = \mathbf{G}_{kj} = 0 \end{cases}$$

where all \mathbf{U}_{kj} 's are independently sampled from the uniform distribution $\text{Uniform}(-1, +1)$. Let \mathcal{C}_2 be the rescaling operator that converts a symmetric positive definite

matrix to the corresponding correlation matrix, we further calculate

$$\mathbf{\Sigma} = \mathbf{\Theta} \mathcal{C}_2[(\tilde{\mathbf{G}} + (0.1 - \Lambda_{\min}(\tilde{\mathbf{G}})) \cdot \mathbf{I})^{-1}] \mathbf{\Theta},$$

where $\mathbf{\Theta}$ is the diagonal standard deviation matrix with $\Theta_{jj} = 2^{\frac{2j-d-1}{2(d-1)}}$ for $j = 1, \dots, d$.

We then generate $n = \lceil 14\sqrt{d} \rceil$ independent samples from the t-distribution with 6 degrees of freedom, mean $\mathbf{0}$, and covariance $\mathbf{\Sigma}$. For the EPIC estimator, we set $c = 0.5$ in (13). For the Catoni's M-estimator, we set $K_{\max} = 10$ and $K_{\min} = 0.1$.

B. Timing Performance

We first evaluate the computational performance of the parametric simplex method. For each model, we choose a regularization parameter, which yields approximate $0.05 \cdot d(d-1)$ nonzero off-diagonal entries. The EPIC and CLIME methods are solved by the parametric simplex method, which is described in Appendix B. The GLASSO is solved by the dual block coordinate descent algorithm, which is described in [11]. Table I summarizes the timing performance averaged over 100 replications. To obtain the baseline performance, we solve the CLIME.SC method using the simplex method¹ as suggested in [5]. We see that all four methods greatly outperform the baseline. The EPIC, CLIME.RC, and CLIME.SC methods attain similar timing performance for all settings, and the GLASSO.RC method is more efficient than the others for $d = 201$ and $d = 401$.

C. Parameter Estimation

To select the regularization parameter, we independently generate a validation set of n samples from the same distribution. We tune λ over a refined grid, then the selected optimal regularization parameter is $\hat{\lambda} = \arg\min_{\lambda} \|\hat{\mathbf{\Omega}}^{\lambda} \hat{\mathbf{\Sigma}} - \mathbf{I}\|_{\max}$, where $\hat{\mathbf{\Omega}}^{\lambda}$ denotes the estimated precision matrix of the training set using the regularization parameter λ , and $\hat{\mathbf{\Sigma}}$ denotes the estimated covariance matrix of the validation set using either (1) or (17).

¹The implementation of the simplex method is based on the R packages `linprog` and `lpSolve`.

TABLE II
QUANTITIVE COMPARISON OF DIFFERENT ESTIMATORS ON THE BAND, ERDÖS-RÉNYI, AND SCALE-FREE MODELS.
THE EPIC ESTIMATOR OUTPERFORMS THE COMPETITORS IN ALL SETTINGS

Spectral Norm: $\ \hat{\Omega} - \Omega\ _2$					
Model	d	EPIC	GLASSO.RC	CLIME.RC	CLIME.SC
Band	101	3.3748(0.2081)	4.4360(9.1445)	3.3961(0.4403)	3.6885(0.5850)
	201	3.3283(0.1114)	4.8616(0.0644)	3.4559(0.0979)	4.4789(0.3399)
	401	3.5933(0.5192)	5.1667(0.0354)	4.0623(0.2397)	5.7164(0.9666)
Erdős-Rényi	101	2.1849(0.2281)	2.6681(0.1293)	2.6787(0.8414)	2.3391(0.2976)
	201	1.8322(0.0769)	2.3753(0.0949)	2.0106(0.3943)	2.0528(0.1548)
	401	1.3322(0.1294)	2.4265(0.0564)	2.0051(0.4144)	4.0667(1.1174)
Scale-free	101	2.1113(0.3081)	2.9979(0.1654)	2.0401(0.3703)	2.6541(0.5882)
	201	2.3519(0.1779)	3.2394(0.1078)	2.3785(0.4186)	2.5789(0.5139)
	401	3.2273(0.1201)	4.0105(0.5812)	3.3139(0.5812)	3.9287(1.1750)

TABLE III
QUANTITIVE COMPARISON OF DIFFERENT ESTIMATORS ON THE BAND, ERDÖS-RÉNYI, AND SCALE-FREE MODELS.
THE EPIC ESTIMATOR OUTPERFORMS THE COMPETITORS IN ALL SETTINGS

Frobenius Norm: $\ \hat{\Omega} - \Omega\ _F$					
Model	d	EPIC	GLASSO.RC	CLIME.RC	CLIME.SC
Band	101	9.4307(0.3245)	11.069(0.2618)	9.7538(0.3949)	11.392(0.8319)
	201	12.720(0.2282)	16.135(0.1399)	13.533(0.1898)	14.850(0.6167)
	401	18.298(1.0537)	23.177(0.1957)	20.412(0.2366)	25.254(1.0002)
Erdős-Rényi	101	6.0660(0.1552)	6.8777(0.2115)	6.7097(0.3672)	7.3789(0.4390)
	201	6.7794(0.1632)	8.1531(0.1828)	7.6175(0.2616)	8.3555(0.2844)
	401	7.3497(0.1743)	10.795(0.1323)	8.3869(0.4755)	11.104(0.6069)
Scale-free	101	4.6695(0.2435)	5.6689(0.2344)	4.9658(0.1762)	6.2264(0.3841)
	201	5.6732(0.1782)	7.2768(0.0940)	6.2343(0.2401)	7.2842(0.3310)
	401	7.2979(0.1094)	9.0940(0.0935)	7.3765(0.2328)	9.5396(0.5636)

Tables II and III summarize the numerical results averaged over 100 replications. We see that the EPIC estimator outperforms the GLASSO.RC and CLIME.RC estimators in all settings.

D. Model Selection

To evaluate the model selection performance, we calculate the ROC curve of each obtained regularization path using the false positive rate (FPR) and true positive rate (FNR) defined as follows,

$$\text{F.P.R.} = \frac{\sum_{k,j} I(\hat{\Omega}_{kj}^{(\lambda)} \neq 0, \Omega_{kj} = 0)}{\sum_{k,j} I(\Omega_{kj} = 0)},$$

$$\text{T.P.R.} = \frac{\sum_{k,j} I(\hat{\Omega}_{kj}^{(\lambda)} \neq 0, \Omega_{kj} \neq 0)}{\sum_{k,j} I(\hat{\Omega}_{kj}^{(\lambda)} \neq 0)}.$$

Figure 1 summarizes ROC curves of all methods averaged over 100 replications.² We see that the EPIC estimator outperforms the competing estimators throughout all settings. Similarly,

²The ROC curves from different replications are first aligned by regularization parameters. The averaged ROC curve shows the false positive and true positive rate averaged over all replications w.r.t. each regularization parameter

our method outperforms the sample covariance matrix-based CLIME estimator.

VI. REAL DATA EXAMPLE

To illustrate the effectiveness of the proposed EPIC method, we adopt the sonar dataset from UCI Machine Learning Repository³ [13]. The dataset contains 101 patterns obtained by bouncing sonar signals off a metal cylinder at various angles and under various conditions, and 97 patterns obtained from rocks under similar conditions. Each pattern is a set of 60 features. Each feature represents the logarithm of the energy integrated over a certain period of time within a particular frequency band. Our goal is to discriminate between sonar signals bounced off a metal cylinder and those bounced off a roughly cylindrical rock.

We randomly split the data into two sets. The training set contains 80 metal and 77 rock patterns. The testing set contains 21 metal and 20 rock patterns. Let $\mu^{(k)}$ be the class conditional means of the data where $k = 1$ represents the metal category and $k = 0$ represents the rock category. [5] assume that two classes share the same covariance matrix, and then adopt the sample mean for estimating μ_k 's and the sample covariance matrix-based CLIME estimator for estimating Ω .

³Available at <http://archive.ics.uci.edu/ml/datasets.html>.

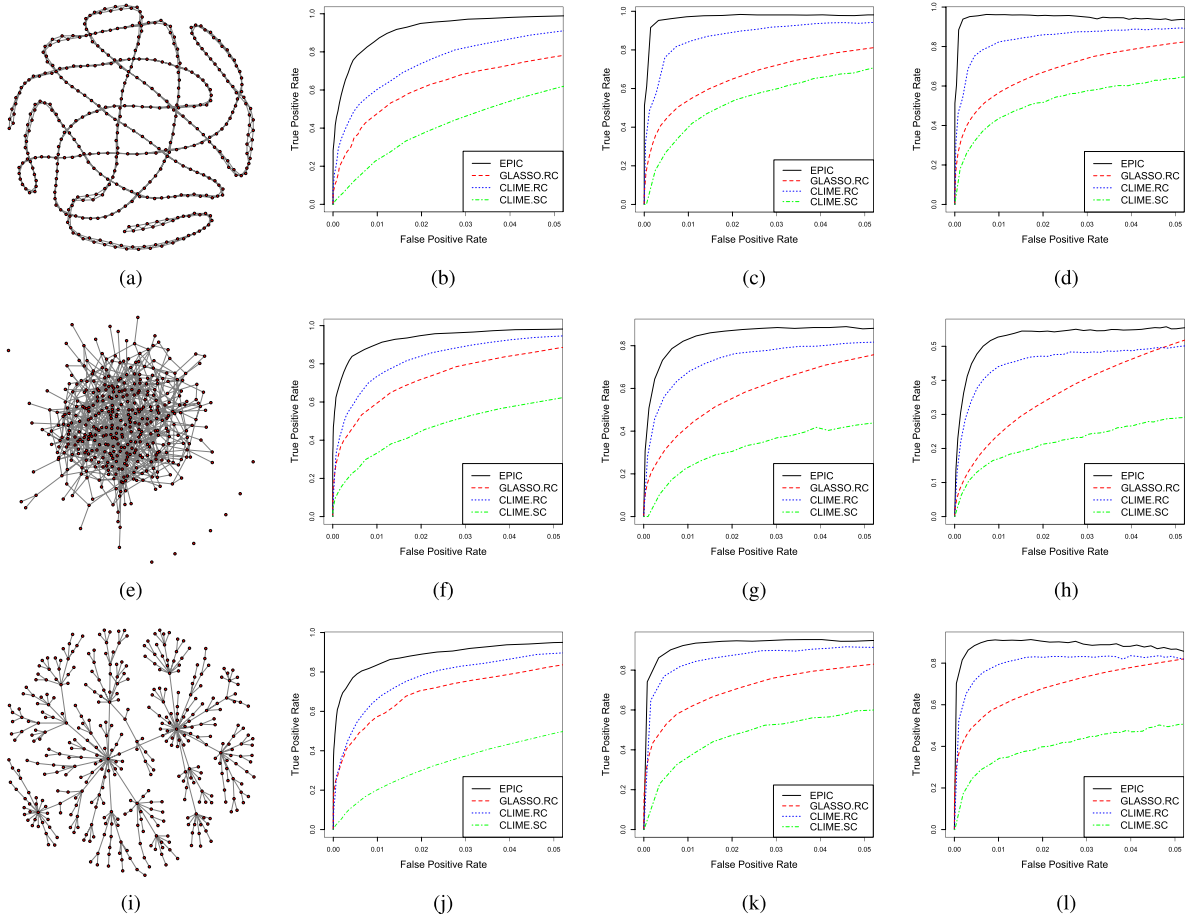


Fig. 1. Three different graph patterns and corresponding average ROC curves. EPIC outperforms the competitors throughout all settings. (a) Band ($d = 401$). (b) Band ($d = 101$). (c) Band ($d = 201$). (d) Band ($d = 401$). (e) Erdős-Rényi ($d = 401$). (f) Erdős-Rényi ($d = 101$). (g) Erdős-Rényi ($d = 201$). (h) Erdős-Rényi ($d = 401$). (i) Scale-free ($d = 401$). (j) Scale-free ($d = 101$). (k) Scale-free ($d = 201$). (l) Scale-free ($d = 401$).

In contrast, we adopt the Catoni's M-estimator for estimating μ_k 's and the EPIC estimator for estimating Ω . We classify a sample x to the metal category if

$$\left(x - \frac{\hat{\mu}^{(1)} + \hat{\mu}^{(0)}}{2}\right)^T \hat{\Omega} (\hat{\mu}^{(1)} - \hat{\mu}^{(0)}) \geq 0,$$

and to the rock category otherwise. We use the testing set to evaluate the performance of the EPIC estimator. For tuning parameter selection, we use a 5-fold cross validation on the training set to pick the regularization parameter λ .

To evaluate the classification performance, we use the criteria of misclassification rate, specificity, sensitivity, and Mathews Correlation Coefficient (MCC). More specifically, let y_i 's and \hat{y}_i 's be true labels and predicted labels of the testing samples, we define

$$\begin{aligned} \text{Misclassification Rate} &= \frac{TP + TN}{TN + TP + FN + FP}, \\ \text{Specificity} &= \frac{TN}{TN + FP}, \quad \text{Sensitivity} = \frac{TP}{TP + FN}, \\ \text{MCC} &= \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}, \end{aligned}$$

where

$$\begin{aligned} TP &= \sum_i I(\hat{y}_i = y_i = 1), \quad FP = \sum_i I(\hat{y}_i = 1, y_i = 0), \\ TN &= \sum_i I(\hat{y}_i = y_i = 0), \quad FN = \sum_i I(\hat{y}_i = 0, y_i = 1). \end{aligned}$$

Table IV summarizes the performance of both methods averaged over 100 replications (with standard errors in parentheses). We see that the EPIC estimator significantly outperforms the competitor on the sensitivity and misclassification rate, but slightly worse on the specificity. The overall classification performance measured by MCC shows that the EPIC estimator has about 8% improvement over the competitor.

VII. DISCUSSION AND CONCLUSION

In this paper, we propose a new sparse precision matrix estimation method for the elliptical family. Our method handles heavy-tailness, and conducts parameter estimation under a calibration framework. We show that the proposed method achieves improved rates of convergence and better finite sample performance than existing methods. The effectiveness of the proposed method is further illustrated by numerical experiments on both simulated and real datasets.

TABLE IV
QUANTITATIVE COMPARISON OF THE EPIC AND SAMPLE COVARIANCE MATRIX-BASED
CLIME ESTIMATORS IN THE SONAR DATA CLASSIFICATION

Method	Misclassification Rate	Specificity	Sensitivity	MCC
EPIC	0.1990(0.0285)	0.7288(0.0499)	0.8579(0.0301)	0.6023(0.0665)
CLIME.SC	0.2362(0.0317)	0.7460(0.0403)	0.7791(0.0429)	0.5288(0.0631)

[25] proposed another calibrated graph estimation method named TIGER for Gaussian family. However, unlike the EPIC estimator, the TIGER method can not handle the elliptical family due to two reasons: (1) The transformed Kendall's tau estimator cannot guarantee the positive semidefiniteness. If we directly plug it into the TIGER method, it makes the TIGER formulation nonconvex. Existing algorithms may not obtain a global solution in polynomial time. (2) The theoretical analysis in [25] is only applicable to the Gaussian family. Theoretical properties of the TIGER method for the elliptical family is unclear.

Another closely related method is the rank-based CLIME method for estimating inverse correlation matrix estimation for the elliptical family [24]. The rank-based CLIME method is based on the formulation in (3) and cannot calibrate the regularization. Furthermore, the rank-based CLIME method can only estimate the inverse correlation matrix. Thus for applications such as the linear discriminant analysis (as is demonstrated in §6) which requires the input to be a precision matrix [2], [30], [35], the rank-based CLIME method is not applicable.

APPENDIX A PROOF OF PROPOSITION III.1

Proof: To show the equivalence between (14) and (15), we only need to verify that the optimal solution $(\hat{\Gamma}_{*j}, \hat{\tau}_j)$ to (15) satisfies

$$\|\hat{\mathbf{Z}}\hat{\Gamma}_{*j} - \mathbf{I}_{*j}\|_{\infty} = \frac{c}{v}\hat{\tau}_j. \quad (\text{A.1})$$

We then prove (A.1) by contradiction. Assuming that there exists some $\bar{\tau}_j \geq 0$ such that

$$\|\hat{\mathbf{Z}}\hat{\Gamma}_{*j} - \mathbf{I}_{*j}\|_{\infty} = \frac{c}{v}\bar{\tau}_j < \frac{c}{v}\hat{\tau}_j, \quad (\text{A.2})$$

(A.2) implies that $(\hat{\Gamma}_{*j}, \bar{\tau}_j)$ is also a feasible solution to (15) and

$$\|\hat{\Gamma}_{*j}\|_1 + c\bar{\tau}_j < \|\hat{\Gamma}_{*j}\|_1 + c\hat{\tau}_j. \quad (\text{A.3})$$

(A.3) contradicts with the fact that $(\hat{\Gamma}_{*j}, \hat{\tau}_j)$ minimizes (15). Thus (A.1) must hold, and (15) is equivalent to (14). \square

APPENDIX B PARAMETRIC SIMPLEX METHOD

We provide a brief description of the parametric simplex method only for self-containedness. More details of the derivation can be found in [34]. We consider the following generic form of linear program,

$$\max_{\mathbf{x} \in \mathbb{R}^m} \mathbf{c}^T \mathbf{x} \quad \text{s.t. } \mathbf{A}\mathbf{x} \leq \mathbf{b}, \mathbf{x} \geq \mathbf{0}, \quad (\text{B.1})$$

where $\mathbf{c} \in \mathbb{R}^m$, $\mathbf{A} \in \mathbb{R}^{n \times m}$, and $\mathbf{b} \in \mathbb{R}^n$. It is well known that (B.1) has a dual formulation as follows,

$$\min_{\mathbf{y} \in \mathbb{R}^n} \mathbf{b}^T \mathbf{y} \quad \text{s.t. } \mathbf{A}^T \mathbf{y} \geq \mathbf{c}, \mathbf{y} \geq \mathbf{0}, \quad (\text{B.2})$$

where $\mathbf{y} = (y_1, \dots, y_n)^T \in \mathbb{R}^n$ are dual variables. The simplex method usually solves either (B.1) or (B.2). It contains two phases: Phase I is to find a feasible initial solution for Phase II; Phase II is an iterative procedure to recover the optimal solution based on the given initial solution.

Different from the simplex method, the parametric simplex method adds some perturbation to (B.1) and (B.2) such that the optimal solutions can be trivially obtained. More specifically, the parametric simplex method solves the following pair of linear programs

$$\max_{\mathbf{x} \in \mathbb{R}^m} (\mathbf{c} + \beta \mathbf{q})^T \mathbf{x} \quad \text{s.t. } \mathbf{A}\mathbf{x} \leq \mathbf{b} + \beta \mathbf{p}, \mathbf{x} \geq \mathbf{0}, \quad (\text{B.3})$$

$$\min_{\mathbf{y} \in \mathbb{R}^n} (\mathbf{b} + \beta \mathbf{p})^T \mathbf{y} \quad \text{s.t. } \mathbf{A}^T \mathbf{y} \geq \mathbf{c} + \beta \mathbf{q}, \mathbf{y} \geq \mathbf{0}, \quad (\text{B.4})$$

where $\beta \geq 0$ is a perturbation parameter, $\mathbf{p} \in \mathbb{R}^n$ and $\mathbf{q} \in \mathbb{R}^m$ are perturbation vectors. When β , \mathbf{p} , and \mathbf{q} are suitably chosen such that $\mathbf{b} + \beta \mathbf{p} \geq \mathbf{0}$ and $\mathbf{c} + \beta \mathbf{q} \leq \mathbf{0}$, $\mathbf{x} = \mathbf{0}$ and $\mathbf{y} = \mathbf{0}$ are the optimal solutions to (B.3) and (B.4) respectively. The parametric simplex method is an iterative procedure, which gradually reduces β to 0 (corresponding to no perturbation) and eventually recovers the optimal solution to (B.1).

To derive the iterative procedure, we first add slack variables $\mathbf{w} = (w_1, \dots, w_n)^T \in \mathbb{R}^n$, and rewrite (B.3) as

$$\max_{\tilde{\mathbf{x}} \in \mathbb{R}^{m+n}} (\tilde{\mathbf{c}} + \beta \tilde{\mathbf{q}})^T \tilde{\mathbf{x}} \quad \text{s.t. } \mathbf{H}\tilde{\mathbf{x}} = \mathbf{b} + \beta \mathbf{p}, \tilde{\mathbf{x}} \geq \mathbf{0}. \quad (\text{B.5})$$

where $\mathbf{H} = [\mathbf{A} \quad \mathbf{I}]$, $\tilde{\mathbf{c}} = (\mathbf{c}^T, \mathbf{0}^T)^T$, $\tilde{\mathbf{q}} = (\mathbf{q}^T, \mathbf{0}^T)^T$ and

$$\begin{aligned} \tilde{\mathbf{x}} &= (\tilde{x}_1, \dots, \tilde{x}_m, \tilde{x}_{m+1}, \dots, \tilde{x}_{m+n})^T \\ &= (x_1, \dots, x_m, w_1, \dots, w_n)^T \in \mathbb{R}^{m+n}. \end{aligned}$$

Since $\mathbf{b} + \beta \mathbf{p} \geq \mathbf{0}$ and $\mathbf{c} + \beta \mathbf{q} \leq \mathbf{0}$, $\tilde{\mathbf{x}} = (\mathbf{0}, \mathbf{b} + \beta \mathbf{p})^T$ is the optimal solution to (B.5). We then divide all variables in $\tilde{\mathbf{x}}$ into a nonbasic group \mathcal{N} and a basic group \mathcal{B} . In particular, $\tilde{x}_1, \dots, \tilde{x}_m$ belong to the nonbasic group denoted by $\tilde{\mathbf{x}}_{\mathcal{N}}$, and $\tilde{x}_{m+1}, \dots, \tilde{x}_{m+n}$ belong to the basic group denoted by $\tilde{\mathbf{x}}_{\mathcal{B}}$. We also divide \mathbf{H} into two submatrices $\mathbf{H}_{\mathcal{N}}$ and $\mathbf{H}_{\mathcal{B}}$, where $\mathbf{H}_{\mathcal{N}}$ contains all columns of \mathbf{H} corresponding to $\tilde{\mathbf{x}}_{\mathcal{N}}$, and $\mathbf{H}_{\mathcal{B}}$ contains all columns of \mathbf{H} corresponding to $\tilde{\mathbf{x}}_{\mathcal{B}}$. We then rewrite the constraint in (B.5) as $\mathbf{H}_{\mathcal{N}}\tilde{\mathbf{x}}_{\mathcal{N}} + \mathbf{H}_{\mathcal{B}}\tilde{\mathbf{x}}_{\mathcal{B}} = \mathbf{b} + \beta \mathbf{p}$. Consequently, we obtain the primal dictionary associated with the basic group \mathcal{B} by

$$\tilde{\mathbf{x}}_{\mathcal{B}} = \tilde{\mathbf{x}}_{\mathcal{B}}^* + \beta \tilde{\mathbf{x}}_{\mathcal{B}} - \mathbf{H}_{\mathcal{B}}^{-1} \mathbf{H}_{\mathcal{N}}, \quad (\text{B.6})$$

$$\phi_P = \phi^* - (\tilde{\mathbf{z}}^* + \beta \tilde{\mathbf{z}})^T \tilde{\mathbf{x}}_{\mathcal{N}}, \quad (\text{B.7})$$

where $\tilde{\mathbf{x}}_{\mathcal{B}}^* = \mathbf{H}_{\mathcal{B}}^{-1}\mathbf{b}$, $\bar{\mathbf{x}}_{\mathcal{B}} = \mathbf{H}_{\mathcal{B}}^{-1}\mathbf{p}$, $\phi_P^* = \tilde{\mathbf{c}}_{\mathcal{B}}^T \mathbf{H}_{\mathcal{B}} \mathbf{b}$, $\tilde{\mathbf{z}}^* = (\mathbf{H}_{\mathcal{B}}^{-1} \mathbf{H}_{\mathcal{N}})^T \tilde{\mathbf{c}}_{\mathcal{B}} - \tilde{\mathbf{c}}_{\mathcal{N}}$, $\tilde{\mathbf{z}}_{\mathcal{N}} = (\mathbf{H}_{\mathcal{B}}^{-1} \mathbf{H}_{\mathcal{N}})^T \tilde{\mathbf{q}}_{\mathcal{B}} - \tilde{\mathbf{q}}_{\mathcal{N}}$, and ϕ_P is the objective value of (B.5) at current iteration.

We then add slack variables $\mathbf{z} = (z_1, \dots, z_m)^T$, and rewrite (16) as

$$\begin{aligned} & \min_{\mathbf{y} \in \mathbb{R}^n} (\tilde{\mathbf{b}} + \beta \tilde{\mathbf{p}})^T \mathbf{y} \\ & \text{s.t. } \mathbf{A} - \mathbf{z}^T \tilde{\mathbf{y}} \geq \mathbf{c} + \beta \mathbf{q}, \quad \tilde{\mathbf{y}} \geq \mathbf{0}. \end{aligned} \quad (\text{B.8})$$

To make the notation consistent with the primal problem, we define

$$\begin{aligned} \tilde{\mathbf{z}} &= (\tilde{z}_1, \dots, \tilde{z}_m, \tilde{z}_{m+1}, \dots, \tilde{z}_{m+n}) \\ &= (z_1, \dots, z_m, y_1, \dots, y_n)^T \in \mathbb{R}^{m+n}. \end{aligned}$$

Similarly we can obtain the dual dictionary associated with the nonbasic variable \mathcal{N} by

$$\tilde{\mathbf{z}}_{\mathcal{N}} = (\tilde{\mathbf{z}}_{\mathcal{N}}^* + \beta \tilde{\mathbf{z}}_{\mathcal{N}}) + (\mathbf{H}_{\mathcal{B}}^{-1} \mathbf{H}_{\mathcal{N}})^T \tilde{\mathbf{z}}_{\mathcal{B}}, \quad (\text{B.9})$$

$$-\phi_D = -\phi_D^* - (\mathbf{x}_{\mathcal{B}}^* + \beta \bar{\mathbf{x}}_{\mathcal{B}})^T \tilde{\mathbf{z}}_{\mathcal{B}}, \quad (\text{B.10})$$

where $\phi_D^* = \tilde{\mathbf{c}}_{\mathcal{B}}^T \mathbf{H}_{\mathcal{B}}^{-1} \mathbf{b}$, and ϕ_D is the objective value of the dual problem at current iteration.

Once we obtain (B.6), (B.7), (B.9), and (B.10), we start to decrease β , and the smallest value of β at current iteration is obtained by

$$\beta^* = \min\{\beta \mid \tilde{\mathbf{x}}_{\mathcal{B}}^* + \beta \bar{\mathbf{x}}_{\mathcal{B}} \geq 0, \tilde{\mathbf{z}}_{\mathcal{N}}^* + \beta \tilde{\mathbf{z}}_{\mathcal{N}} \geq 0\}.$$

we then swap a pair of basic and nonbasic variables in \mathcal{B} and \mathcal{N} and update the primal and dual dictionaries such that β can be decreased to β^* . See more details on updating the dictionaries in [34]. By repeating the above procedure, we eventually decrease β to 0. The parametric simplex method guarantees the feasibility and optimality for both (B.3) and (B.4) in each iteration, and eventually obtain the optimal solution to the original problem (B.3).

Since the parametric simplex method starts with all zero solutions, it can recover the optimal solution only in a few iterations when the optimal solution is very sparse. That naturally fits into the sparse estimation problems such as the EPIC method. Moreover, if we rewrite (16) in the same form as (B.3), we need to set $\mathbf{p} = (\mathbf{0}^T, \mathbf{e}^T, 0)^T$ and start with $\beta = 1$. Since $\mathbf{c} = (-\mathbf{1}^T, -\mathbf{c})^T$, we can set $\mathbf{q} = \mathbf{0}$, i.e., we do not need perturbation on \mathbf{c} . Thus the computation in each iteration can be further simplified due to the sparsity of \mathbf{p} and \mathbf{q} .

Remark B.1: For sparse estimation problems, Phase I of the simplex method does not guarantee the sparseness of the initial solution. As a result, Phase II may start with a dense initial solution, and gradually reduce the sparsity of the solution. Thus the overall convergence of the simplex method often requires a large number of iterations when the optimal solution is very sparse.

APPENDIX C

SMOOTHED PROXIMAL GRADIENT ALGORITHM

We first apply the smoothing approach in [28] to obtain a smooth surrogate of the matrix ℓ_1 norm based on the Fenchel dual representation,

$$\|\mathbf{\Omega} - \tilde{\mathbf{\Omega}}\|_{\eta} = \min_{\|\mathbf{U}\|_{\infty} \leq 1} \text{tr}(\mathbf{U}^T (\mathbf{\Omega} - \tilde{\mathbf{\Omega}})) + \frac{\eta}{2} \|\mathbf{U}\|_{\text{F}}^2, \quad (\text{C.1})$$

where $\eta > 0$ is a smoothing parameter. (C.1) has a closed form solution $\hat{\mathbf{U}}$ as follows,

$$\hat{\mathbf{U}}_{kj}^{\mathbf{\Omega}} = \text{sign}(\tilde{\mathbf{U}}_{kj}^{\mathbf{\Omega}}) \cdot \max\{|\tilde{\mathbf{U}}_{kj}^{\mathbf{\Omega}}| - \gamma_k\}. \quad (\text{C.2})$$

where $\tilde{\mathbf{U}}^{\mathbf{\Omega}} = (\mathbf{\Omega} - \tilde{\mathbf{\Omega}})/\eta$, and γ_k is the minimum positive value such that $\|\hat{\mathbf{U}}\|_{\infty} = \max_k \|\hat{\mathbf{U}}_{k*}\|_1 \leq 1$. See [9] for an efficient algorithms to find γ_k with the average computational complexity of $O(d^2)$. As is shown in [28], the smooth surrogate $\|\mathbf{\Omega} - \tilde{\mathbf{\Omega}}\|_{\eta}$ is smooth, convex, and has a simple form gradient as

$$\mathbf{G}(\mathbf{\Omega}) = \frac{\partial \|\mathbf{\Omega} - \tilde{\mathbf{\Omega}}\|_{\eta}}{\partial \mathbf{\Omega}} = \hat{\mathbf{U}}_{kj}^{\mathbf{\Omega}}.$$

Since $\hat{\mathbf{U}}_{kj}^{\mathbf{\Omega}}$ is obtained by the soft-thresholding in (C.2), we have $\mathbf{G}(\mathbf{\Omega})$ continuous in $\mathbf{\Omega}$ with the Lipschitz constant η^{-1} . Motivated by these good computational properties, we consider the following optimization problem instead of (18),

$$\bar{\mathbf{\Omega}} = \underset{\mathbf{\Omega} = \mathbf{\Omega}^T}{\text{argmin}} \|\mathbf{\Omega} - \tilde{\mathbf{\Omega}}\|_{\eta}. \quad (\text{C.3})$$

To solve (C.3), we adopt the accelerated projected gradient algorithm proposed in [27]. More specifically, we define two sequences of auxiliary variables $\{\mathbf{M}^{(t)}\}$ and $\{\mathbf{W}^{(t)}\}$ with $\mathbf{M}^{(0)} = \mathbf{W}^{(0)} = \mathbf{\Omega}^{(0)}$, and a sequence of weights $\{\theta_t = 2/(1+t)\}$. For the t^{th} iteration, we first calculate the auxiliary variable $\mathbf{M}^{(t)}$ as

$$\mathbf{M}^{(t)} = (1 - \theta_t) \mathbf{\Omega}^{(t-1)} + \theta_t \mathbf{W}^{(t-1)}.$$

We then calculate the auxiliary variable $\mathbf{W}^{(t)}$ as

$$\begin{aligned} \mathbf{W}^{(t)} &= \underset{\mathbf{W} = \mathbf{W}^T}{\text{argmin}} \|\mathbf{W}^{(t-1)} - \tilde{\mathbf{\Omega}}\|_{\eta} + \text{tr}((\mathbf{W} - \mathbf{W}^{(t-1)})^T \mathbf{G}(\mathbf{M}^{(t)})) \\ &\quad + \frac{1}{2\eta_t \theta_t} \|\mathbf{W} - \mathbf{W}^{(t-1)}\|_{\text{F}}^2 = \frac{1}{2} \left(\mathbf{W}^{(t-1)} - \frac{\eta_t}{\theta_t} \mathbf{G}(\mathbf{M}^{(t)}) \right) \\ &\quad + \frac{1}{2} \left(\mathbf{W}^{(t-1)} - \frac{\eta_t}{\theta_t} \mathbf{G}(\mathbf{M}^{(t)}) \right)^T, \end{aligned}$$

where η_t is the step size. We can either choose $\eta_t = \eta$ in all iterations or estimate η_t 's by the back-tracking line search for better empirical performance [4]. At last, we calculate $\mathbf{\Omega}^{(t)}$ as,

$$\mathbf{\Omega}^{(t)} = (1 - \theta_t) \mathbf{\Omega}^{(t-1)} + \theta_t \mathbf{W}^{(t)}.$$

The next theorem provides the convergence rate of the algorithm with respect to minimizing (18).

Theorem C.1: Given the desired accuracy ε such that $\|\mathbf{\Omega}^{(t)} - \tilde{\mathbf{\Omega}}\|_1 - \|\hat{\mathbf{\Omega}} - \tilde{\mathbf{\Omega}}\|_1 < \varepsilon$, let $\eta = d^{-1}\varepsilon/2$, we need the number of iterations to be at most

$$t = 2\sqrt{2d} \|\mathbf{\Omega}^{(0)} - \bar{\mathbf{\Omega}}\|_{\text{F}} \cdot \varepsilon^{-1} - 1 = O(\varepsilon^{-1}).$$

Proof: Due to the fact that $\|\mathbf{A}\|_{\text{F}} \leq d \|\mathbf{A}\|_{\infty}$, a direct consequence of (C.1) is the following uniform bound

$$\|\mathbf{\Omega} - \tilde{\mathbf{\Omega}}\|_1 - d\eta \leq \|\mathbf{\Omega} - \tilde{\mathbf{\Omega}}\|_{\eta} \leq \|\mathbf{\Omega} - \tilde{\mathbf{\Omega}}\|_1.$$

Then we consider the following decomposition

$$\begin{aligned}
& \|\boldsymbol{\Omega}^{(t)} - \tilde{\boldsymbol{\Omega}}\|_1 - \|\hat{\boldsymbol{\Omega}} - \tilde{\boldsymbol{\Omega}}\|_1 \\
&= \|\boldsymbol{\Omega}^{(t)} - \tilde{\boldsymbol{\Omega}}\|_1 - \|\tilde{\boldsymbol{\Omega}} - \tilde{\boldsymbol{\Omega}}\|_\eta + \|\tilde{\boldsymbol{\Omega}} - \tilde{\boldsymbol{\Omega}}\|_\eta - \|\hat{\boldsymbol{\Omega}} - \tilde{\boldsymbol{\Omega}}\|_1 \\
&\leq \|\boldsymbol{\Omega}^{(t)} - \tilde{\boldsymbol{\Omega}}\|_\eta - \|\tilde{\boldsymbol{\Omega}} - \tilde{\boldsymbol{\Omega}}\|_\eta + d\eta \\
&\leq \frac{2\|\boldsymbol{\Omega}^{(t)} - \tilde{\boldsymbol{\Omega}}\|_F^2}{(t+1)^2\eta} + d\eta,
\end{aligned}$$

where the last inequality comes from the result established in [27],

$$\|\boldsymbol{\Omega}^{(t)} - \tilde{\boldsymbol{\Omega}}\|_\eta - \|\tilde{\boldsymbol{\Omega}} - \tilde{\boldsymbol{\Omega}}\|_\eta \leq \frac{2\|\boldsymbol{\Omega}^{(0)} - \tilde{\boldsymbol{\Omega}}\|_F^2}{(t+1)^2\eta}.$$

Thus given $d\eta = \epsilon/2$, we only need

$$\frac{4d\|\boldsymbol{\Omega}^{(0)} - \tilde{\boldsymbol{\Omega}}\|_F^2}{(t+1)^2} \leq \frac{\epsilon^2}{2}. \quad (\text{C.4})$$

By solving (C.4), we obtain

$$t \leq \frac{2\sqrt{2d}\|\boldsymbol{\Omega}^{(0)} - \tilde{\boldsymbol{\Omega}}\|_F}{\epsilon} - 1. \quad \square$$

Theorem C.1 guarantees that the above algorithm achieves the optimal rate of convergence for minimizing (18) over the class of all first-order computational algorithms.

APPENDIX D PROOF OF LEMMA 1

Proof: [7] shows that there exist universal constants κ_3 and κ_4 such that

$$\mathbb{P}(|\hat{\mu}_j - \mu_j| \leq \kappa_3\theta_{\max}\epsilon) \geq 1 - \exp(-n\epsilon^2), \quad (\text{D.1})$$

$$\mathbb{P}(|\hat{m}_j - \mathbb{E}X_j^2| \leq \kappa_4\sqrt{K}\epsilon) \geq 1 - \exp(-n\epsilon^2). \quad (\text{D.2})$$

We then define the following events

$$\begin{aligned}
\mathcal{C}_1 &= \{|\hat{\mu}_j - \mu_j| \leq \mu_{\max}\}, \\
\mathcal{C}_2 &= \{|\hat{\mu}_j - \mu_j| \leq \kappa_3\theta_{\max}\epsilon\}, \\
\mathcal{C}_3 &= \{|\hat{m}_j - \mathbb{E}X_j^2| \leq \kappa_4\sqrt{K}\epsilon\}, \\
\mathcal{C}_4 &= \{|\hat{\theta}_j - \theta_j| \leq \theta_{\min}\}.
\end{aligned}$$

Conditioning on \mathcal{C}_1 , we have

$$\begin{aligned}
|\hat{\mu}_j^2 - \mu_j^2| &= |\hat{\mu}_j - \mu_j| \cdot |\hat{\mu}_j + \mu_j| \\
&\leq |\hat{\mu}_j - \mu_j| \cdot |\hat{\mu}_j - \mu_j| + |\hat{\mu}_j - \mu_j| \cdot |2\mu_j| \\
&\leq (2\mu_{\max} + |\hat{\mu}_j - \mu_j|)|\hat{\mu}_j - \mu_j| \\
&\leq 3\mu_{\max}|\hat{\mu}_j - \mu_j|.
\end{aligned} \quad (\text{D.3})$$

Conditioning on \mathcal{C}_2 and \mathcal{C}_3 , (D.3) implies

$$\begin{aligned}
|\hat{\theta}_j^2 - \theta_j^2| &= |\hat{m}_j - \hat{\mu}_j^2 - \mathbb{E}X_j^2 + \mu_j^2| \\
&\leq |\hat{m}_j - \mathbb{E}X_j^2| + |\hat{\mu}_j^2 - \mu_j^2| \\
&\leq (3\mu_{\max}\kappa_3\theta_{\max} + \kappa_4\sqrt{K})\epsilon.
\end{aligned} \quad (\text{D.4})$$

(D.4) further implies

$$\begin{aligned}
|\hat{\theta}_j - \theta_j| &\leq \frac{|\hat{\theta}_j^2 - \theta_j^2|}{\hat{\theta}_j + \theta_j} \leq \frac{|\hat{\theta}_j^2 - \theta_j^2|}{\theta_j} \\
&\leq \frac{(3\mu_{\max}\kappa_3\theta_{\max} + \kappa_4\sqrt{K})\epsilon}{\theta_{\min}}.
\end{aligned} \quad (\text{D.5})$$

Conditioning \mathcal{C}_4 , (D.5) implies

$$\begin{aligned}
|\hat{\theta}_j^{-1} - \theta_j^{-1}| &= \frac{|\hat{\theta}_j - \theta_j|}{\hat{\theta}_j\theta_j} = \frac{|\hat{\theta}_j - \theta_j|}{(\hat{\theta}_j - \theta_j)\theta_j + 2\theta_j^2} \\
&\leq \frac{(3\mu_{\max}\kappa_3\theta_{\max} + \kappa_4\sqrt{K})\epsilon}{(\hat{\theta}_j - \theta_j)\theta_j\theta_{\min} + 2\theta_j^2\theta_{\min}} \\
&\leq \frac{(3\mu_{\max}\kappa_3\theta_{\max} + \kappa_4\sqrt{K})\epsilon}{2\theta_j^2\theta_{\min}} \\
&\leq \frac{(3\mu_{\max}\kappa_3\theta_{\max} + \kappa_4\sqrt{K})\epsilon}{2\theta_{\min}^3}.
\end{aligned} \quad (\text{D.6})$$

Combining (D.1), (D.2), and (D.6), for small enough ϵ such that

$$\epsilon \leq \min \left\{ \frac{\mu_{\max}}{\kappa_3\theta_{\max}}, \frac{\theta_{\min}^2}{3\mu_{\max}\kappa_3\theta_{\max} + \kappa_4\sqrt{K}} \right\}, \quad (\text{D.7})$$

we have

$$\mathbb{P}\left(|\hat{\theta}_j^{-1} - \theta_j^{-1}| \leq \frac{(3\mu_{\max}\kappa_3\theta_{\max} + \kappa_4\sqrt{K})\epsilon}{\theta_{\min}^3}\right) \geq 1 - 2\exp(-4n\epsilon^2). \quad (\text{D.8})$$

By taking the union bound of (D.8), we have

$$\begin{aligned}
\mathbb{P}\left(\max_{1 \leq j \leq d} |\hat{\theta}_j^{-1} - \theta_j^{-1}| \leq \frac{(3\mu_{\max}\kappa_3\theta_{\max} + \kappa_4\sqrt{K})\epsilon}{\theta_{\min}^3}\right) \\
\geq 1 - 2\exp(-4n\epsilon^2 + \log d).
\end{aligned}$$

If we take $\epsilon = \sqrt{\log d/n}$, then (D.7) implies that we need n large enough such that

$$n \geq \max \left\{ \frac{\kappa_3^2\theta_{\max}^2}{\mu_{\max}^2}, \frac{(3\mu_{\max}\kappa_3\theta_{\max} + \kappa_4\sqrt{K})^2}{\theta_{\min}^4} \right\} \cdot \log d.$$

Taking $\kappa_2 = (3\mu_{\max}\kappa_3\theta_{\max} + \kappa_4\sqrt{K})/\theta_{\min}^3$, we then have

$$\mathbb{P}\left(\max_{1 \leq j \leq d} |\hat{\theta}_j^{-1} - \theta_j^{-1}| \leq \kappa_2\sqrt{\frac{\log d}{n}}\right) \geq 1 - \frac{2}{d}.$$

(19) is a direct result in [24], therefore its proof is omitted. \square

APPENDIX E PROOF OF THEOREM IV.1

Proof: We first define the following pair of orthogonal subspaces $(\mathcal{S}_j, \mathcal{S}_j^\perp)$,

$$\begin{aligned}
\mathcal{S}_j &= \{\mathbf{v} \in \mathbb{R}^d \mid v_k = 0 \text{ for all } \mathbf{\Gamma}_{kj} = 0\}, \\
\mathcal{S}_j^\perp &= \{\mathbf{v} \in \mathbb{R}^d \mid v_k = 0 \text{ for all } \mathbf{\Gamma}_{kj} \neq 0\}.
\end{aligned}$$

We will use $(\mathcal{S}_j, \mathcal{S}_j^\perp)$ to exploit the sparseness of $\mathbf{\Gamma}_{*j}$. We then define the following event

$$\mathcal{D}_1 = \{\|\hat{\mathbf{Z}} - \mathbf{Z}\|_{\max} \leq \lambda\}.$$

Conditioning on \mathcal{D}_1 , we have

$$\begin{aligned}
\|\hat{\mathbf{Z}}\mathbf{\Gamma}_{*j} - \mathbf{I}_{*j}\|_\infty &= \|(\hat{\mathbf{Z}} - \mathbf{Z})\mathbf{\Gamma}_{*j}\|_\infty \\
&\leq \|\mathbf{\Gamma}_{*j}\|_1 \|\hat{\mathbf{Z}} - \mathbf{Z}\|_{\max} \leq \lambda \|\mathbf{\Gamma}_{*j}\|_1.
\end{aligned} \quad (\text{E.1})$$

Now let $\tau_j = \|\mathbf{\Gamma}_{*j}\|_1$, (E.1) implies that $(\mathbf{\Gamma}_{*j}, \tau_j)$ is a feasible solution to (13). Since $(\widehat{\mathbf{\Gamma}}_{*j}, \widehat{\tau}_j)$ is the empirical minimizer, we have

$$\begin{aligned} \|\widehat{\mathbf{\Gamma}}_{S_{jj}}\|_1 + \|\widehat{\mathbf{\Gamma}}_{S_{jj}^\perp}\|_1 + c\widehat{\tau}_j &= \|\widehat{\mathbf{\Gamma}}_{*j}\|_1 + c\widehat{\tau}_j \\ &\leq \|\mathbf{\Gamma}_{*j}\|_1 + c\tau_j = \|\mathbf{\Gamma}_{S_{jj}}\|_1 + c\tau_j, \end{aligned} \quad (\text{E.2})$$

where the last equality comes from the fact that $\mathbf{\Gamma}_{S_{jj}^\perp} = \mathbf{0}$. Let $\widehat{\mathbf{\Delta}} = \widehat{\mathbf{\Gamma}} - \mathbf{\Gamma}$ be the estimation error, (E.2) implies

$$\begin{aligned} \|\widehat{\mathbf{\Gamma}}_{S_{jj}^\perp}\|_1 &\leq \|\mathbf{\Gamma}_{S_{jj}}\|_1 - \|\widehat{\mathbf{\Gamma}}_{S_{jj}}\|_1 + c(\tau_j - \widehat{\tau}_j) \\ &\leq \|\widehat{\mathbf{\Delta}}_{S_{jj}}\|_1 + c(\tau_j - \widehat{\tau}_j) \\ &\leq \|\widehat{\mathbf{\Delta}}_{S_{jj}}\|_1 + c(\|\mathbf{\Gamma}_{*j}\|_1 - \|\widehat{\mathbf{\Gamma}}_{*j}\|_1) \\ &\stackrel{(i)}{\leq} \|\widehat{\mathbf{\Delta}}_{S_{jj}}\|_1 + c\|\widehat{\mathbf{\Delta}}_{*j}\|_1 \\ &\stackrel{(ii)}{\leq} (1+c)\|\widehat{\mathbf{\Delta}}_{S_{jj}}\|_1 + c\|\widehat{\mathbf{\Delta}}_{S_{jj}^\perp}\|_1, \end{aligned} \quad (\text{E.3})$$

where (i) comes from the constraint in (13): $\|\widehat{\mathbf{\Gamma}}_{*j}\|_1 \leq \widehat{\tau}_j$, and (ii) comes from the fact $\|\widehat{\mathbf{\Delta}}_{*j}\|_1 = \|\widehat{\mathbf{\Delta}}_{S_{jj}}\|_1 + \|\widehat{\mathbf{\Delta}}_{S_{jj}^\perp}\|_1$. Combining the fact $\|\widehat{\mathbf{\Delta}}_{S_{jj}^\perp}\|_1 = \|\widehat{\mathbf{\Gamma}}_{S_{jj}^\perp} - \mathbf{0}\|_1 = \|\widehat{\mathbf{\Gamma}}_{S_{jj}^\perp}\|_1$ with (E.3), we have

$$\|\widehat{\mathbf{\Delta}}_{S_{jj}^\perp}\|_1 \leq \bar{c}\|\widehat{\mathbf{\Delta}}_{S_{jj}}\|_1, \quad (\text{E.4})$$

where $\bar{c} = (1+c)/(1-c)$. (E.4) implies that $\widehat{\mathbf{\Delta}}_{*j}$ belongs the following cone shape set

$$\mathcal{M}_j^{\bar{c}} = \left\{ \mathbf{v} \in \mathbb{R}^d \setminus \{\mathbf{0}\} \mid \|\mathbf{v}_{S_{jj}^\perp}\|_1 \leq \bar{c}\|\mathbf{v}_{S_{jj}}\|_1 \right\}.$$

The following lemma characterizes an important property of $\mathcal{M}_j^{\bar{c}}$ when \mathcal{D}_1 holds.

Lemma E.1: Suppose that $\mathbf{X} \sim EC(\boldsymbol{\mu}, \boldsymbol{\xi}, \boldsymbol{\Sigma})$, and (A.1) and \mathcal{D}_1 hold. Given any $\mathbf{v} \in \mathcal{M}_j^{\bar{c}}$, for small enough λ such that $2(1+\bar{c})^2 s \lambda \kappa_u \leq 1$, we have

$$\min_{\mathbf{v} \in \mathcal{M}_j^{\bar{c}}} \mathbf{v}^T \widehat{\mathbf{Z}} \mathbf{v} \geq \frac{\|\mathbf{v}\|_2^2}{2\kappa_u}. \quad (\text{E.5})$$

The proof of Lemma E.1 is provided in Appendix E.1. Since $\widehat{\mathbf{\Delta}}_{*j}$ exactly belongs to $\mathcal{M}_j^{\bar{c}}$, we have a simple variant of (E.5) as

$$\begin{aligned} \|\widehat{\mathbf{\Delta}}_{*j}\|_1 \|\widehat{\mathbf{Z}} \widehat{\mathbf{\Delta}}_{*j}\|_\infty &\geq \widehat{\mathbf{\Delta}}_{*j}^T \widehat{\mathbf{Z}} \widehat{\mathbf{\Delta}}_{*j} \geq \frac{\|\widehat{\mathbf{\Delta}}_{*j}\|_2^2}{2\kappa_u} \\ &\geq \frac{\|\widehat{\mathbf{\Delta}}_{S_{jj}}\|_2^2}{2\kappa_u} \geq \frac{\|\widehat{\mathbf{\Delta}}_{S_{jj}}\|_1^2}{2s\kappa_u}, \end{aligned} \quad (\text{E.6})$$

where the last inequality comes from the fact that $\widehat{\mathbf{\Delta}}_{S_{jj}}$ has at most s nonzero entries. Since

$$\begin{aligned} \|\widehat{\mathbf{Z}} \widehat{\mathbf{\Delta}}_{*j}\|_\infty &\leq \|\widehat{\mathbf{Z}} \widehat{\mathbf{\Gamma}}_{*j} - \mathbf{I}_{*j}\|_\infty + \|\widehat{\mathbf{Z}} \mathbf{\Gamma}_{*j} - \mathbf{I}_{*j}\|_\infty \\ &\leq \lambda(\widehat{\tau}_j + \tau_j) \\ &\leq \lambda(2\widehat{\tau}_j + \tau_j - \widehat{\tau}_j) \\ &\leq \lambda(2\widehat{\tau}_j + \|\widehat{\mathbf{\Delta}}_{*j}\|_1) \\ &\leq \lambda(2\widehat{\tau}_j + (1+\bar{c})\|\widehat{\mathbf{\Delta}}_{S_{jj}}\|_1), \end{aligned} \quad (\text{E.7})$$

where the last inequality comes from (E.4). Combining (E.6) and (E.7), we have

$$\|\widehat{\mathbf{Z}} \widehat{\mathbf{\Delta}}_{*j}\|_\infty \leq \lambda(2\widehat{\tau}_j + 2(1+\bar{c})\kappa_u s \|\widehat{\mathbf{Z}} \widehat{\mathbf{\Delta}}_{*j}\|_\infty). \quad (\text{E.8})$$

Assuming that $1 - 2(1+\bar{c})\kappa_u s \lambda = \delta_1 > 0$, (E.8) implies

$$\|\widehat{\mathbf{Z}} \widehat{\mathbf{\Delta}}_{*j}\|_\infty \leq 2\delta_1^{-1} \lambda \widehat{\tau}_j. \quad (\text{E.9})$$

Combining (E.6) and (E.9), we have

$$\begin{aligned} c\widehat{\tau}_j &\leq \|\widehat{\mathbf{\Delta}}_{S_{jj}}\|_1 + c\tau_j \leq 2\kappa_u s \|\widehat{\mathbf{Z}} \widehat{\mathbf{\Delta}}_{*j}\|_\infty + c\tau_j \\ &\leq 4\kappa_u s \delta_1^{-1} \lambda \widehat{\tau}_j + c\tau_j. \end{aligned} \quad (\text{E.10})$$

Assuming that $1 - 4\kappa_u s \delta_1^{-1} c^{-1} \lambda = \delta_2 > 0$, (E.10) implies

$$\widehat{\tau}_j \leq \delta_2^{-1} \tau_j. \quad (\text{E.11})$$

Recall $\lambda = \kappa_1 \sqrt{\log d/n}$, in order to secure

$$\begin{aligned} 1 - 2(1+\bar{c})\kappa_u s \lambda &= \delta_1 > 0, \\ c - 4\kappa_u s \delta_1^{-1} \lambda &= \delta_2 > 0, \\ 2(1+\bar{c})^2 s \lambda \kappa_u &\leq 1, \end{aligned}$$

we need large enough n such that

$$n \geq \max \left\{ 4(1-\delta_1)^{-2} (1+\bar{c})^2 \kappa_u^2, 16(c-\delta_2)^{-2} \kappa_u^2 \delta_1^{-2}, 4(1+\bar{c})^4 \kappa_u \kappa_1 \right\} \cdot s^2 \log d.$$

Combining (E.9) and (E.11), we have

$$\|\widehat{\mathbf{Z}} \widehat{\mathbf{\Delta}}_{*j}\|_\infty \leq 2\delta_1^{-1} \delta_2^{-1} \lambda \tau_j. \quad (\text{E.12})$$

Combining (E.4), (E.6), and (E.12), we obtain

$$\|\widehat{\mathbf{\Delta}}_{*j}\|_1 \leq 4(1+\bar{c})\kappa_u \delta_1^{-1} \delta_2^{-1} s \lambda \tau_j. \quad (\text{E.13})$$

Combining (E.12) and (E.13), we have

$$\begin{aligned} \widehat{\mathbf{\Delta}}_{*j}^T \widehat{\mathbf{Z}} \widehat{\mathbf{\Delta}}_{*j} &\leq \|\widehat{\mathbf{\Delta}}_{*j}\|_1 \cdot \|\widehat{\mathbf{Z}} \widehat{\mathbf{\Delta}}_{*j}\|_\infty \\ &\leq 8(1+\bar{c})\kappa_u \delta_1^{-2} \delta_2^{-2} s \lambda^2 \tau_j^2. \end{aligned} \quad (\text{E.14})$$

By Lemma E.1 again, (E.14) implies

$$\|\widehat{\mathbf{\Delta}}_{*j}\|_2^2 \leq 16(1+\bar{c})\kappa_u^2 \delta_1^{-2} \delta_2^{-2} s \lambda^2 \tau_j^2. \quad (\text{E.15})$$

Let $\kappa_5 = 4(1+\bar{c})\kappa_u \delta_1^{-1} \delta_2^{-1} \kappa_1$ and $\kappa_6 = 16(1+\bar{c})\kappa_u^2 \delta_1^{-2} \delta_2^{-2} \kappa_1^2$. Recall $\lambda = \kappa_1 \sqrt{\log d/n}$, by definition of the matrix ℓ_1 and Frobenius norms, (E.13) and (E.15) imply

$$\begin{aligned} \|\widehat{\mathbf{\Delta}}\|_1 &= \max_j \|\widehat{\mathbf{\Delta}}_{*j}\|_1 \leq 4(1+\bar{c})\kappa_u \delta_1^{-1} \delta_2^{-1} s \lambda \max_j \tau_j \\ &\leq \kappa_5 \cdot M \cdot s \sqrt{\frac{\log d}{n}}, \end{aligned} \quad (\text{E.16})$$

and

$$\begin{aligned} \frac{1}{d} \|\widehat{\mathbf{\Delta}}\|_F^2 &= \frac{1}{d} \sum_j \|\widehat{\mathbf{\Delta}}_{*j}\|_2^2 \leq 16(1+\bar{c})\kappa_u^2 \delta_1^{-2} \delta_2^{-2} s \lambda^2 \tau_j^2 \\ &\leq \kappa_6 \cdot M^2 \cdot \frac{s \log d}{n}. \end{aligned} \quad (\text{E.17})$$

Now we start to derive the error bound of $\tilde{\mathbf{\Omega}}$ obtained by the ensemble rule. We have the following decomposition

$$\begin{aligned}\tilde{\mathbf{\Omega}} - \mathbf{\Omega} &= \hat{\mathbf{\Theta}}^{-1} \hat{\mathbf{\Gamma}} \hat{\mathbf{\Theta}}^{-1} - \mathbf{\Theta}^{-1} \mathbf{\Gamma} \mathbf{\Theta}^{-1} \\ &= (\hat{\mathbf{\Theta}}^{-1} - \mathbf{\Theta}^{-1} + \mathbf{\Theta}^{-1})(\hat{\mathbf{\Gamma}} - \mathbf{\Gamma} + \mathbf{\Gamma}) \\ &\quad \cdot (\hat{\mathbf{\Theta}}^{-1} - \mathbf{\Theta}^{-1} + \mathbf{\Theta}^{-1}) - \mathbf{\Theta}^{-1} \mathbf{\Gamma} \mathbf{\Theta}^{-1} \\ &= (\hat{\mathbf{\Theta}}^{-1} - \mathbf{\Theta}^{-1})(\hat{\mathbf{\Gamma}} - \mathbf{\Gamma})(\hat{\mathbf{\Theta}}^{-1} - \mathbf{\Theta}^{-1}) \\ &\quad + (\hat{\mathbf{\Theta}}^{-1} - \mathbf{\Theta}^{-1})(\hat{\mathbf{\Gamma}} - \mathbf{\Gamma})\mathbf{\Theta}^{-1} + (\hat{\mathbf{\Theta}}^{-1} - \mathbf{\Theta}^{-1})\mathbf{\Gamma}\mathbf{\Theta}^{-1} \\ &\quad + (\hat{\mathbf{\Theta}}^{-1} - \mathbf{\Theta}^{-1})\mathbf{\Gamma}(\hat{\mathbf{\Theta}}^{-1} - \mathbf{\Theta}^{-1}) + \mathbf{\Theta}^{-1}(\hat{\mathbf{\Gamma}} - \mathbf{\Gamma})\mathbf{\Theta}^{-1} \\ &\quad + \mathbf{\Theta}^{-1}(\hat{\mathbf{\Gamma}} - \mathbf{\Gamma})(\hat{\mathbf{\Theta}}^{-1} - \mathbf{\Theta}^{-1}) \\ &\quad + \mathbf{\Theta}^{-1}\mathbf{\Gamma}(\hat{\mathbf{\Theta}}^{-1} - \mathbf{\Theta}^{-1}).\end{aligned}\quad (\text{E.18})$$

Moreover, for any $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{R}^{d \times d}$, where \mathbf{A} and \mathbf{C} are diagonal matrices, we have

$$\|\mathbf{ABC}\|_1 \leq \|\mathbf{A}\|_{\max} \cdot \|\mathbf{B}\|_1 \cdot \|\mathbf{C}\|_{\max}, \quad (\text{E.19})$$

$$\|\mathbf{ABC}\|_F \leq \|\mathbf{A}\|_{\max} \cdot \|\mathbf{B}\|_F \cdot \|\mathbf{C}\|_{\max}. \quad (\text{E.20})$$

Here we define the following event

$$\mathcal{D}_2 = \left\{ \|\hat{\mathbf{\Theta}}^{-1} - \mathbf{\Theta}^{-1}\|_{\max} \leq \kappa_2 \sqrt{\frac{\log d}{n}} \right\}.$$

Thus conditioning \mathcal{D}_2 , (E.16), (E.18), and (E.19) imply

$$\begin{aligned}\|\tilde{\mathbf{\Omega}} - \mathbf{\Omega}\|_1 &\leq \kappa_2^2 \kappa_5 \cdot \frac{\log d}{n} \cdot M \cdot s \sqrt{\frac{\log d}{n}} \\ &\quad + \frac{\kappa_2 \kappa_5}{\theta_{\min}} \cdot \sqrt{\frac{\log d}{n}} \cdot M \cdot s \sqrt{\frac{\log d}{n}} \\ &\quad + \kappa_2^2 \cdot M \cdot \frac{\log d}{n} + \frac{\kappa_2}{\theta_{\min}} \cdot M \cdot \sqrt{\frac{\log d}{n}} \\ &\quad + \frac{\kappa_2 \kappa_5}{\theta_{\min}} \cdot \sqrt{\frac{\log d}{n}} \cdot M \cdot s \sqrt{\frac{\log d}{n}} \\ &\quad + \frac{\kappa_5}{\theta_{\min}^2} M \cdot s \sqrt{\frac{\log d}{n}} + \frac{\kappa_2}{\theta_{\min}} \cdot M \cdot \sqrt{\frac{\log d}{n}}.\end{aligned}\quad (\text{E.21})$$

If (A.4): $s^2 \log d / n \rightarrow 0$ holds, then (E.21) is determined by the slowest rate $M s \sqrt{\log d / n}$. Thus for large enough n , there exists a universal constant C_4 such that

$$\|\tilde{\mathbf{\Omega}} - \mathbf{\Omega}\|_1 \leq C_4 \cdot M \cdot s \sqrt{\frac{\log d}{n}}. \quad (\text{E.22})$$

Similarly, conditioning on \mathcal{D}_2 , (E.17), (E.18), (E.20) and the fact $\|\mathbf{\Gamma}\|_F \leq M \sqrt{d}$ imply

$$\begin{aligned}\|\tilde{\mathbf{\Omega}} - \mathbf{\Omega}\|_F &\leq \kappa_2^2 \kappa_6 \cdot \frac{\log d}{n} \cdot M \cdot \sqrt{\frac{ds \log d}{n}} \\ &\quad + \frac{\kappa_2 \kappa_6}{\theta_{\min}} \cdot \sqrt{\frac{\log d}{n}} \cdot M \cdot \sqrt{\frac{ds \log d}{n}} \\ &\quad + \kappa_2^2 \cdot M \sqrt{d} \cdot \frac{\log d}{n} + \frac{\kappa_2}{\theta_{\min}} \cdot M \sqrt{d} \cdot \sqrt{\frac{\log d}{n}} \\ &\quad + \frac{\kappa_2 \kappa_6}{\theta_{\min}} \cdot \sqrt{\frac{\log d}{n}} \cdot M \cdot \sqrt{\frac{ds \log d}{n}} \\ &\quad + \frac{\kappa_6}{\theta_{\min}^2} M \cdot \sqrt{\frac{ds \log d}{n}} + \frac{\kappa_2}{\theta_{\min}} \cdot M \sqrt{d} \cdot \sqrt{\frac{\log d}{n}}.\end{aligned}\quad (\text{E.23})$$

Again if (A.4) holds, then (E.23) is determined by the slowest rate $M \sqrt{ds \log d / n}$. Thus for large enough n , there exists a universal constant C_2 such that

$$\frac{1}{d} \|\tilde{\mathbf{\Omega}} - \mathbf{\Omega}\|_F^2 \leq C_2 \cdot M^2 \cdot \frac{s \log d}{n}. \quad (\text{E.24})$$

We then proceed to prove the error bound of $\hat{\mathbf{\Omega}}$ obtained by the symmetrization procedure (18). Let $C_1 = 2C_4$, if we choose the matrix ℓ_1 norm as $\|\cdot\|_*$ in (18), we have

$$\begin{aligned}\|\hat{\mathbf{\Omega}} - \mathbf{\Omega}\|_1 &\leq \|\tilde{\mathbf{\Omega}} - \hat{\mathbf{\Omega}}\|_1 + \|\tilde{\mathbf{\Omega}} - \mathbf{\Omega}\|_1 \\ &\leq 2\|\tilde{\mathbf{\Omega}} - \mathbf{\Omega}\|_1 \leq C_1 \cdot M \cdot s \sqrt{\frac{\log d}{n}},\end{aligned}\quad (\text{E.25})$$

where the second inequality comes from the fact that $\mathbf{\Omega}$ is a feasible solution to (18), and $\hat{\mathbf{\Omega}}$ is the empirical minimizer. If we choose the Frobenius norm as $\|\cdot\|_*$ in (18), using the fact that the Frobenius norm projection is contractive, we have

$$\frac{1}{d} \|\hat{\mathbf{\Omega}} - \mathbf{\Omega}\|_F^2 \leq \frac{1}{d} \|\tilde{\mathbf{\Omega}} - \mathbf{\Omega}\|_F^2 \leq C_2 M^2 \frac{s \log d}{n}. \quad (\text{E.26})$$

All above analysis are conditioned on \mathcal{D}_1 and \mathcal{D}_2 . Thus combining Lemma 1 with (E.25) and (E.26), we have

$$\mathbb{P}\left(\|\hat{\mathbf{\Omega}} - \mathbf{\Omega}\|_p \leq C_1 M s \sqrt{\frac{\log d}{n}}\right) \geq 1 - \frac{3}{d}, \quad (\text{E.27})$$

$$\mathbb{P}\left(\frac{1}{d} \|\hat{\mathbf{\Omega}} - \mathbf{\Omega}\|_F^2 \leq C_2 M^2 \frac{s \log d}{n}\right) \geq 1 - \frac{3}{d}. \quad (\text{E.28})$$

where $p = 1, 2$, and (E.27) comes from the fact that $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_1$ for any symmetric matrix \mathbf{A} . \square

APPENDIX F

PROOF OF LEMMA E.1

Proof: Since $\Lambda_{\min}(\mathbf{Z}) = 1/\Lambda_{\max}(\mathbf{\Gamma}) \geq \kappa_u^{-1}$, we have

$$\begin{aligned}\mathbf{v}^T \hat{\mathbf{Z}} \mathbf{v} &= \mathbf{v}^T \mathbf{Z} \mathbf{v} - \mathbf{v}^T (\mathbf{Z} - \hat{\mathbf{Z}}) \mathbf{v} \\ &\geq \kappa_u^{-1} \|\mathbf{v}\|_2^2 - \|\mathbf{v}\|_1^2 \cdot \|\mathbf{Z} - \hat{\mathbf{Z}}\|_{\max}.\end{aligned}\quad (\text{F.1})$$

Since $\mathbf{v} \in \mathcal{M}_j^c$, we have $\|\mathbf{v}_{\mathcal{S}_j^c}\|_1 \leq \bar{c} \|\mathbf{v}_{\mathcal{S}_j}\|_1$, which implies

$$\begin{aligned}\|\mathbf{v}\|_1 &= \|\mathbf{v}_{\mathcal{S}_j}\|_1 + \|\mathbf{v}_{\mathcal{S}_j^c}\|_1 \leq (1 + \bar{c}) \|\mathbf{v}_{\mathcal{S}_j}\|_1 \\ &\leq (1 + \bar{c}) \sqrt{s} \|\mathbf{v}_{\mathcal{S}_j}\|_2,\end{aligned}\quad (\text{F.2})$$

where the last inequality comes from the fact that there are at most s nonzero entries in $\mathbf{v}_{\mathcal{S}_j}$. Then combining (F.1) and (F.2), we have

$$\begin{aligned}\mathbf{v}^T \hat{\mathbf{Z}} \mathbf{v} &\geq \kappa_u^{-1} \|\mathbf{v}\|_2^2 - (1 + \bar{c})^2 \|\mathbf{v}_{\mathcal{S}_j}\|_1^2 \cdot \|\mathbf{Z} - \hat{\mathbf{Z}}\|_{\max} \\ &\geq \kappa_u^{-1} \|\mathbf{v}\|_2^2 - (1 + \bar{c})^2 s \lambda \|\mathbf{v}_{\mathcal{S}_j}\|_2^2.\end{aligned}\quad (\text{F.3})$$

Since we have $2(1 + \bar{c})^2 s \lambda \kappa_u \leq 1$, (F.3) implies

$$\mathbf{v}^T \hat{\mathbf{Z}} \mathbf{v} \geq \frac{\|\mathbf{v}\|_2^2}{2\kappa_u}. \quad \square$$

APPENDIX G

PROOF OF THEOREM IV.2

Proof: Our following analysis also assumes that $\mathcal{D}_1 = \{\|\hat{\mathbf{Z}} - \mathbf{Z}\|_{\max} \leq \lambda\}$ holds. Since \mathcal{D}_1 implies (E.1),

$$\|\hat{\mathbf{Z}} \mathbf{\Gamma}_{*j} - \mathbf{I}_{*j}\|_{\infty} \leq \lambda \tau_j, \quad \forall j = 1, \dots, d,$$

where $\tau_j = \|\mathbf{\Gamma}_{*j}\|_1$. Then $(\mathbf{\Gamma}_{*j}, \tau_j)$ is a feasible solution to (13), which implies

$$\begin{aligned} \|\widehat{\mathbf{Z}}(\widehat{\mathbf{\Gamma}}_{*j} - \mathbf{\Gamma}_{*j})\|_\infty &\leq \|\widehat{\mathbf{Z}}\widehat{\mathbf{\Gamma}}_{*j} - \mathbf{I}_{*j}\|_\infty + \|\widehat{\mathbf{Z}}\mathbf{\Gamma}_{*j} - \mathbf{I}_{*j}\|_\infty \\ &\leq \lambda\widehat{\tau}_j + \lambda\tau_j. \end{aligned} \quad (\text{G.1})$$

Moreover, we have

$$\begin{aligned} (1+c)\|\widehat{\mathbf{\Gamma}}_{*j}\|_1 &\leq \|\widehat{\mathbf{\Gamma}}_{*j}\|_1 + c\widehat{\tau}_j \leq \|\mathbf{\Gamma}_{*j}\|_1 + c\tau_j \\ &= (1+c)\|\mathbf{\Gamma}_{*j}\|_1 = (1+c)\tau_j, \end{aligned}$$

which further implies

$$\|\widehat{\mathbf{\Gamma}}_{*j}\|_1 \leq \|\mathbf{\Gamma}_{*j}\|_1, \quad \|\widehat{\mathbf{\Gamma}}_{*j} - \mathbf{\Gamma}_{*j}\|_1 \leq 2\|\mathbf{\Gamma}_{*j}\|_1. \quad (\text{G.2})$$

Combining (G.1) and (G.2), we have

$$\begin{aligned} \|\mathbf{Z}(\widehat{\mathbf{\Gamma}}_{*j} - \mathbf{\Gamma}_{*j})\|_\infty &\leq \|\widehat{\mathbf{Z}}(\widehat{\mathbf{\Gamma}}_{*j} - \mathbf{\Gamma}_{*j})\|_\infty + \|(\widehat{\mathbf{Z}} - \mathbf{Z})(\widehat{\mathbf{\Gamma}}_{*j} - \mathbf{\Gamma}_{*j})\|_\infty \\ &\leq \lambda\tau + \lambda\widehat{\tau}_j + \|\widehat{\mathbf{Z}} - \mathbf{Z}\|_{\max} \|\widehat{\mathbf{\Gamma}}_{*j} - \mathbf{\Gamma}_{*j}\|_1 \\ &\leq 3\lambda\tau + \lambda\widehat{\tau}_j \leq \frac{(1+4c)\lambda\tau_j}{c}, \end{aligned} \quad (\text{G.3})$$

where the last inequality comes from

$$c\widehat{\tau}_j \leq \|\widehat{\mathbf{\Gamma}}_{*j}\|_1 + c\widehat{\tau}_j \leq \|\mathbf{\Gamma}_{*j}\|_1 + c\tau_j = (1+c)\tau_j.$$

By (G.3), we have

$$\begin{aligned} \|\widehat{\mathbf{\Gamma}}_{*j} - \mathbf{\Gamma}_{*j}\|_\infty &\leq \|\mathbf{\Gamma}_{*j}\|_1 \|\mathbf{Z}(\widehat{\mathbf{\Gamma}}_{*j} - \mathbf{\Gamma}_{*j})\|_\infty \\ &\leq \frac{\lambda(1+4c)\tau_j}{c} \cdot \|\mathbf{\Gamma}_{*j}\|_1 \leq \frac{\lambda(1+4c)\tau_j^2}{c}. \end{aligned} \quad (\text{G.4})$$

Recall $\lambda = \kappa_1\sqrt{\log d/n}$, by the definition of the max norm and (G.4), we have

$$\|\widehat{\mathbf{\Gamma}} - \mathbf{\Gamma}\|_{\max} \leq \kappa_6 \cdot M^2 \cdot \sqrt{\frac{\log d}{n}}, \quad (\text{G.5})$$

where $\kappa_7 = \kappa_1(1+4c)/c$. Since for any $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{R}^{d \times d}$, where \mathbf{A} and \mathbf{C} are diagonal matrices, we have

$$\|\mathbf{ABC}\|_{\max} \leq \|\mathbf{A}\|_{\max} \cdot \|\mathbf{B}\|_{\max} \cdot \|\mathbf{C}\|_{\max}. \quad (\text{G.6})$$

Conditioning on

$$\mathcal{D}_2 = \left\{ \|\widehat{\mathbf{\Theta}}^{-1} - \mathbf{\Theta}^{-1}\|_{\max} \leq \kappa_2 \sqrt{\frac{\log d}{n}} \right\}, \quad (\text{G.7})$$

(G.6), (E.18) and the fact $\|\mathbf{\Gamma}\|_{\max} \leq M$ imply

$$\begin{aligned} \|\widetilde{\mathbf{\Omega}} - \mathbf{\Omega}\|_{\max} &\leq \kappa_2^2 \kappa_7 \cdot \frac{\log d}{n} \cdot M^2 \cdot \sqrt{\frac{\log d}{n}} \\ &\quad + \frac{\kappa_2 \kappa_7}{\theta_{\min}} \cdot \sqrt{\frac{\log d}{n}} \cdot M^2 \cdot \sqrt{\frac{\log d}{n}} \\ &\quad + \kappa_2^2 \cdot M \cdot \frac{\log d}{n} + \frac{\kappa_2}{\theta_{\min}} \cdot M \cdot \sqrt{\frac{\log d}{n}} \\ &\quad + \frac{\kappa_2 \kappa_7}{\theta_{\min}} \cdot \sqrt{\frac{\log d}{n}} \cdot M^2 \cdot \sqrt{\frac{\log d}{n}} \\ &\quad + \frac{\kappa_7}{\theta_{\min}^2} M^2 \cdot \sqrt{\frac{\log d}{n}} + \frac{\kappa_2}{\theta_{\min}} \cdot M \cdot \sqrt{\frac{\log d}{n}}. \end{aligned} \quad (\text{G.8})$$

Again if (A.4): $s^2 \log d/n \rightarrow 0$ holds, then (G.8) is determined by the slowest rate $M^2 \sqrt{\log d/n}$. Thus for large enough n , if we choose the max norm as $\|\cdot\|_*$ in (18), we have

$$\begin{aligned} \|\widehat{\mathbf{\Omega}} - \mathbf{\Omega}\|_{\max} &\leq \|\widetilde{\mathbf{\Omega}} - \widehat{\mathbf{\Omega}}\|_{\max} + \|\widetilde{\mathbf{\Omega}} - \mathbf{\Omega}\|_{\max} \\ &\leq 2\|\widetilde{\mathbf{\Omega}} - \mathbf{\Omega}\|_{\max} \leq 2\kappa_7 \cdot M^2 \cdot \sqrt{\frac{\log d}{n}}, \end{aligned} \quad (\text{G.9})$$

where the second inequality comes from the fact that $\mathbf{\Omega}$ is a feasible solution to (18), and $\widehat{\mathbf{\Omega}}$ is the empirical minimizer.

Note that the results obtained here only depend on \mathcal{D}_1 and \mathcal{D}_2 . Thus by Lemma 1 and (G.9), let $C_3 = 2\kappa_7$, we have

$$\mathbb{P}\left(\|\widehat{\mathbf{\Omega}} - \mathbf{\Omega}\|_{\max} \leq C_3 \cdot M^2 \cdot \sqrt{\frac{\log d}{n}}\right) \geq 1 - \frac{3}{d}.$$

To show the partial consistency in graph estimation $\mathbb{P}(E \subseteq \widehat{E}) \rightarrow 1$, we follow a similar argument to Theorem 4 in [25]. Therefore the proof is omitted. \square

REFERENCES

- [1] O. Banerjee, L. El Ghaoui, and A. d'Aspremont, "Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data," *J. Mach. Learn. Res.*, vol. 9, pp. 485–516, Jun. 2008.
- [2] P. J. Bickel and E. Levina, "Some theory for Fisher's linear discriminant function, 'naive Bayes', and some alternatives when there are many more variables than observations," *Bernoulli*, vol. 10, no. 6, pp. 989–1010, 2004.
- [3] D. M. Blei and J. D. Lafferty, "A correlated topic model of science," *Ann. Appl. Statist.*, vol. 1, no. 1, pp. 17–35, 2007.
- [4] S. Boyd and L. Vandenberghe, *Convex Optimization*, 2nd ed. Cambridge, U.K.: Cambridge Univ. Press, 2009.
- [5] T. Cai, W. Liu, and X. Luo, "A constrained ℓ_1 minimization approach to sparse precision matrix estimation," *J. Amer. Statist. Assoc.*, vol. 106, no. 494, pp. 594–607, 2011.
- [6] S. Cambanis, S. Huang, and G. Simons, "On the theory of elliptically contoured distributions," *J. Multivariate Anal.*, vol. 11, no. 3, pp. 368–385, 1981.
- [7] O. Catoni, "Challenging the empirical mean and empirical variance: A deviation study," *Ann. Inst. Henri Poincaré Probab. Statist.*, vol. 48, no. 4, pp. 1148–1185, 2012.
- [8] A. P. Dempster, "Covariance selection," *Biometrics*, vol. 28, no. 1, pp. 157–175, 1972.
- [9] J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra, "Efficient projections onto the ℓ_1 -ball for learning in high dimensions," in *Proc. 25th Int. Conf. Mach. Learn.*, 2008, pp. 272–279.
- [10] K.-T. Fang, S. Kotz, and K. W. Ng, *Symmetric Multivariate and Related Distributions. Monographs on Statistics and Applied Probability*, vol. 36. London, U.K.: Chapman & Hall, 1990.
- [11] J. Friedman, T. Hastie, and R. Tibshirani, "Sparse inverse covariance estimation with the graphical lasso," *Biostatistics*, vol. 9, no. 3, pp. 432–441, 2008.
- [12] E. Gautier and A. B. Tsybakov, "High-dimensional instrumental variables regression and confidence sets," ENSAE ParisTech, Malakoff, France, Tech. Rep. arxiv.org, 2011.
- [13] R. P. Gorman and T. J. Sejnowski, "Analysis of hidden units in a layered network trained to classify sonar targets," *Neural Netw.*, vol. 1, no. 1, pp. 75–89, 1988.
- [14] A. K. Gupta, T. Varga, and T. Bodnar, *Elliptically Contoured Models in Statistics and Portfolio Theory*. New York, NY, USA: Springer-Verlag, 2013.
- [15] J. Honorio, L. Ortiz, D. Samaras, N. Paragios, and R. Goldstein, "Sparse and locally constant Gaussian graphical models," in *Advances in Neural Information Processing Systems 22*. Red Hook, NY, USA: Curran Associates, 2009.
- [16] C.-J. Hsieh, I. S. Dhillon, P. K. Ravikumar, and M. A. Sustik, "Sparse inverse covariance matrix estimation using quadratic approximation," in *Advances in Neural Information Processing Systems*, vol. 24. Red Hook, NY, USA: Curran Associates, 2011, pp. 2330–2338.

- [17] H. Hult and F. Lindskog, "Multivariate extremes, aggregation and dependence in elliptical distributions," *Adv. Appl. Probab.*, vol. 34, no. 3, pp. 587–608, 2002.
- [18] W. H. Kruskal, "Ordinal measures of association," *J. Amer. Statist. Assoc.*, vol. 53, no. 284, pp. 814–861, 1958.
- [19] W. J. Krzanowski, *Principles of Multivariate Analysis*. Oxford, U.K.: Clarendon, 2000.
- [20] C. Lam and J. Fan, "Sparsistency and rates of convergence in large covariance matrix estimation," *Ann. Statist.*, vol. 37, no. 6B, pp. 4254–4278, 2009.
- [21] S. L. Lauritzen, *Graphical Models*, vol. 17. London, U.K.: Oxford Univ. Press, 1996.
- [22] X. Li, T. Zhao, X. Yuan, and H. Liu, "The flare Package for High-dimensional Sparse Linear Regression in R," *J. Mach. Learn. Res.*, 2014.
- [23] H. Liu, F. Han, M. Yuan, J. Lafferty, and L. Wasserman, "High-dimensional semiparametric Gaussian copula graphical models," *Ann. Statist.*, vol. 40, no. 4, pp. 2293–2326, 2012.
- [24] H. Liu, F. Han, and C.-H. Zhang, "Transelliptical graphical models," in *Advances in Neural Information Processing Systems 25*. Red Hook, NY, USA: Curran Associates, 2012.
- [25] H. Liu and L. Wang, "TIGER: A tuning-insensitive approach for optimally estimating Gaussian graphical models," Massachusetts Inst. Technol., Cambridge, MA, USA, Tech. Rep., 2012.
- [26] H. Liu, L. Wang, and T. Zhao, "Sparse covariance matrix estimation with eigenvalue constraints," *J. Comput. Graph. Statist.*, vol. 23, no. 2, pp. 439–459, 2014.
- [27] Y. E. Nesterov, "An approach to constructing optimal methods for minimization of smooth convex functions," *Ekonomika Matematicheskie Metody*, vol. 24, no. 3, pp. 509–517, 1988.
- [28] Y. Nesterov, "Smooth minimization of non-smooth functions," *Math. Program.*, vol. 103, no. 1, pp. 127–152, 2005.
- [29] A. J. Rothman, P. J. Bickel, E. Levina, and J. Zhu, "Sparse permutation invariant covariance estimation," *Electron. J. Statist.*, vol. 2, pp. 494–515, 2008.
- [30] J. Shao, Y. Wang, X. Deng, and S. Wang, "Sparse linear discriminant analysis by thresholding for high dimensional data," *Ann. Statist.*, vol. 39, no. 2, pp. 1241–1265, 2011.
- [31] J. Stoer, R. Bulirsch, R. Bartels, W. Gautschi, and C. Witzgall, *Introduction to Numerical Analysis*, vol. 2. New York, NY, USA: Springer-Verlag, 1993.
- [32] T. Sun and C. Zhang, "Scaled sparse linear regression," *Biometrika*, vol. 99, no. 4, p. 879, 2012.
- [33] T. Tokuda, B. Goodrich, I. Van Mechelen, A. Gelman, and F. Tuerlinckx, "Visualizing distributions of covariance matrices," Columbia Univ., New York, NY, USA, Tech. Rep., 2011.
- [34] R. J. Vanderbei, *Linear Programming: Foundations and Extensions*. New York, NY, USA: Springer-Verlag, 2008.
- [35] H. Wakaki, "Discriminant analysis under elliptical populations," *Hiroshima Math. J.*, vol. 24, no. 2, pp. 257–298, 1994.
- [36] A. Wille *et al.*, "Sparse graphical Gaussian modeling of the isoprenoid gene network in arabidopsis thaliana," *Genome Biol.*, vol. 25, no. 5, p. R92, 2004.
- [37] M. Yuan, "High dimensional inverse covariance matrix estimation via linear programming," *J. Mach. Learn. Res.*, vol. 11, pp. 2261–2286, Mar. 2010.
- [38] M. Yuan and Y. Lin, "Model selection and estimation in the Gaussian graphical model," *Biometrika*, vol. 94, no. 1, pp. 19–35, 2007.
- [39] T. Zhao and H. Liu, "Sparse inverse covariance estimation with calibration," in *Advances in Neural Information Processing Systems 26*. Red Hook, NY, USA: Curran Associates, 2013.
- [40] T. Zhao, H. Liu, K. Roeder, J. Lafferty, and L. Wasserman, "The huge package for high-dimensional undirected graph estimation in R," *J. Mach. Learn. Res.*, vol. 13, no. 1, pp. 1059–1062, 2012.

Tuo Zhao received his B.S. and M.S. degrees in Computer Science from Harbin Institute of Technology, and his second M.S. degree in Applied Math from University of Minnesota.

He is currently a Ph.D. Candidate in Department of Computer Science at Johns Hopkins University. He is also a visiting student in Department of Operations Research and Financial Engineering at Princeton University. His research focuses on large-scale semiparametric and nonparametric learning and applications to high throughput genomics and neuroimaging.

Han Liu received a joint Ph.D. degree in Machine Learning and Statistics from the Carnegie Mellon University, Pittsburgh, PA, USA in 2011.

He is currently an Assistant Professor of Statistical Machine Learning in the Department of Operations Research and Financial Engineering at Princeton University, Princeton, NJ. He is also an adjunct Professor in the Department of Biostatistics and Department of Computer Science at Johns Hopkins University. He built and is serving as the principal investigator of the Statistical Machine Learning (SMiLe) lab at Princeton University. His research interests include high dimensional semiparametric inference, statistical optimization, Big Data inferential analysis.