

THE EFFECT OF L_1 PENALIZATION ON CONDITION NUMBER CONSTRAINED ESTIMATION OF PRECISION MATRIX

Xiao Guo¹ and Chunming Zhang^{2,3}

¹*University of Science and Technology of China*, ²*Nankai University*
and ³*University of Wisconsin-Madison*

Abstract: Estimation of large precision matrices is fundamental to high-dimensional inference. An important issue is to deal with ill-conditioning of the precision matrix estimate, typically encountered in finite-samples, but rarely studied in the literature. In this paper, we focus on estimating the precision matrix by imposing a bound on the condition number of the estimate, which effectively ensures well-conditioning. Specifically, we propose a correlation-based estimator, constrained with both the condition number and the L_1 penalty, yielding a precision matrix estimator with theoretically guaranteed rate of convergence. This result further enables us to demonstrate that incorporating the L_1 penalty is necessary for achieving consistency of the resulting estimator in typical high-dimensional settings, while inconsistency will occur when the L_1 penalty is absent. An algorithm based on the alternating direction method of multipliers is developed to implement the proposed method, which reveals the satisfactory performance in simulation studies. An application of the method to a call center data is illustrated.

Key words and phrases: Condition number, covariance matrix, penalization, precision matrix, sparsity.

1. Introduction

Estimation of a large precision matrix has been an important and challenging problem with applications in many scientific fields. For example, in linear discriminant analysis, optimal portfolio selection, recovery of the structure of undirected Gaussian graphical model, and detection of activated brain regions for neuroimaging data estimation of the precision matrix is needed. Given n i.i.d. p_n -variate random vectors, the inverse of the sample covariance matrix, \mathbf{S}_n^{-1} , is commonly used for estimating the precision matrix Σ_0^{-1} , where Σ_0 is the true covariance matrix. When the dimension p_n is fixed, \mathbf{S}_n is consistent for Σ_0 , but when $p_n > n$, the singularity of \mathbf{S}_n makes its inverse unavailable. Even if $p_n \leq n$, \mathbf{S}_n^{-1} may be inconsistent when $\lim_{n \rightarrow \infty} p_n/n = c$ for a constant $c \in (0, 1]$

(Marčenko and Pastur (1967)). Besides the potential issue of inconsistency associated with precision matrix estimator, another issue is that its condition number is large, or inflated in practice.

In the literature, several approaches for estimating the covariance matrix have been developed. See Bickel and Levina (2008), Cai and Liu (2011), Cai and Zhou (2012), Liu, Wang and Zhao (2014), Rothman (2012), Rothman, Levina and Zhu (2009), and Xue, Ma and Zou (2012), among others. The inverses of these estimators are consistent for the precision matrix. At the same time, regularization schemes have been proposed to estimate the precision matrix directly. For example, Meinshausen and Bühlmann (2006) proposed an L_1 penalized regression approach, which was extended by Peng et al. (2009). Other works utilizing the penalized log-likelihood approach include Banerjee, El Ghaoui and d'Aspremont (2008), Friedman, Hastie and Tibshirani (2008), and Lam and Fan (2009). These methods simultaneously recover the sparsity structure of the precision matrix. Although the aforementioned estimators are consistent, the problem of ill-conditioning was not taken into consideration. To protect the condition number of the precision matrix estimate from being inflated, a natural way is to shrink the eigenvalues of the estimator. For example, Won et al. (2013) developed an estimator of Σ_0 by imposing a bound on the condition number of the estimator, but did not address the issue of inconsistency.

In this work, we focus on estimating the precision matrix with a condition number constraint. We consider a correlation-based estimator of the precision matrix with the condition number constraint, and study its asymptotic properties. We incorporate the L_1 penalty with the proposed estimator and examine its effect on the consistency of the condition number constrained estimator. We show that if the L_1 penalty is absent, the estimator is consistent only in restrictive cases and inconsistent in many circumstances. Under regularity conditions, we find the convergence rate of our estimator with the L_1 penalty incorporated in high-dimensional cases, allowing $\lim_{n \rightarrow \infty} p_n/n > 0$.

Our estimator with the L_1 penalty is asymptotically equivalent to the correlation-based SPICE estimator developed in Rothman et al. (2008), but has an advantage in that it possesses a constrained condition number and enjoys better finite-sample performance. To implement our estimation method, we develop an algorithm based on the alternating direction method of multipliers (Boyd et al. (2010)). Simulations and data analyses reveal the satisfactory performance of our proposed estimator with the L_1 penalty included.

The rest of this paper is organized as follows. Section 2 proposes the condi-

tion number constrained estimator of the precision matrix, details situations in which consistency and inconsistency take place and derives the convergence rate of the estimator with the L_1 penalty incorporated. Section 3 develops the algorithm for the estimator $\widehat{\Theta}_{\text{prop-2}}$ defined at (2.5) and (2.6). Section 4 discusses data-driven choices of the tuning parameters. Section 5 presents simulations and Section 6 analyzes data. The supplementary material includes the proofs of the results.

We introduce some notation here. For any set G , denote by $|G|$ the cardinality of G . For matrices A and B of size $m \times m$, $\lambda_{\min}(A)$ and $\lambda_{\max}(A)$ denote the smallest and largest eigenvalues of A , respectively, $A \succ 0$ means that A is positive definite, and $A \otimes B$ is the Kronecker matrix product. We write $A(i, j)$ for the element of A in the i th row and j th column. The trace of A is denoted by $\text{tr}(A)$ and $\det(A)$ is the determinant of A . The off-diagonal elementwise L_1 norm of A is

$$|A|_1 = \sum_{1 \leq i \neq j \leq m} |A(i, j)|. \tag{1.1}$$

The L_2 , L_∞ , and Frobenius norms of A are $\|A\|_2 = \{\lambda_{\max}(A^T A)\}^{1/2}$, $\|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^m |A(i, j)|$, and $\|A\|_F = \{\text{tr}(A^T A)\}^{1/2}$, respectively. Denote by \mathbf{I}_m the $m \times m$ identity matrix and by $\mathbf{e}_{q,m}$ the q th column of \mathbf{I}_m . The empirical spectral distribution of A is $\mathbb{F}^A(x) = m^{-1}|\{j \leq m : \lambda_j \leq x\}|$, where $\{\lambda_j\}_{j=1}^m$ are the eigenvalues of A . For a sequence of random distribution functions $\{F_n\}_{n \geq 1}$ and a deterministic distribution function F_0 , we write $F_n(x) \xrightarrow{P} F_0(x)$ as $n \rightarrow \infty$ at any continuous point x of F_0 to denote that F_n converges weakly to F_0 in probability. For a vector $\mathbf{v} = (v_1, \dots, v_m)^T$, the L_1 norm is $\|\mathbf{v}\|_1 = \sum_{i=1}^m |v_i|$. For two sequences of positive real numbers $\{a_n\}_{n=1}^\infty$ and $\{b_n\}_{n=1}^\infty$, $a_n \asymp b_n$ denotes that $a_n = O(b_n)$ and $b_n = O(a_n)$. In the following, C and c are generic finite constants that may vary from place to place and do not depend on n .

2. Condition Number Constrained Estimator of Σ_0^{-1}

Throughout, $\mathbf{X}_1, \dots, \mathbf{X}_n$ are i.i.d. p_n -variate random vectors, with $\mathbf{X}_i = (X_{i,1}, \dots, X_{i,p_n})^T$. We write $\boldsymbol{\mu}_0 = E(\mathbf{X}_i) \in \mathbb{R}^{p_n}$ and $\Sigma_0 = \text{cov}(\mathbf{X}_i, \mathbf{X}_i)$ and assume Σ_0 is positive definite. If $\mathbf{X}_1, \dots, \mathbf{X}_n$ are normally distributed, the log-likelihood is

$$\ell_n(\boldsymbol{\mu}_0, \Sigma_0) = -\frac{np_n}{2} \log(2\pi)$$

$$-\frac{1}{2} \left[-n \log\{\det(\boldsymbol{\Sigma}_0^{-1})\} + \sum_{i=1}^n (\mathbf{X}_i - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}_0^{-1} (\mathbf{X}_i - \boldsymbol{\mu}_0) \right].$$

Write the maximum likelihood estimators (MLEs) of $\boldsymbol{\mu}_0$ and $\boldsymbol{\Sigma}_0$ as $\bar{\mathbf{X}}$ and \mathbf{S}_n . If we replace $\boldsymbol{\mu}_0$ with $\bar{\mathbf{X}}$, the Gaussian log-likelihood function is

$$\ell_n(\bar{\mathbf{X}}, \boldsymbol{\Sigma}_0) = -\frac{np_n}{2} \log(2\pi) - \frac{n}{2} [-\log\{\det(\boldsymbol{\Sigma}_0^{-1})\} + \text{tr}(\mathbf{S}_n \boldsymbol{\Sigma}_0^{-1})].$$

Let $\boldsymbol{\Theta}_0 = \boldsymbol{\Sigma}_0^{-1}$ be the precision matrix of \mathbf{X}_i . We focus on estimating $\boldsymbol{\Theta}_0$. It is known that, when p_n is fixed, \mathbf{S}_n^{-1} is a well-behaved estimator of $\boldsymbol{\Theta}_0$. However, when $\lim_{n \rightarrow \infty} p_n/n = c$ with a constant $c \in (0, 1]$, \mathbf{S}_n^{-1} may be ill-conditioned (Marčenko and Pastur (1967)), i.e. the condition number is inflated where condition number of a positive-definite matrix $\boldsymbol{\Sigma}$ is defined as $\text{cond}(\boldsymbol{\Sigma}) = \lambda_{\max}(\boldsymbol{\Sigma})/\lambda_{\min}(\boldsymbol{\Sigma})$. Particularly, if $p_n > n$, then \mathbf{S}_n is not invertible. When the dimension is large, \mathbf{S}_n^{-1} is either numerically unavailable or ill-conditioned.

We consider estimating $\boldsymbol{\Theta}_0$ by imposing a condition number constraint while minimizing the negative Gaussian log-likelihood function. Although the underlying distribution of $\mathbf{X}_1, \dots, \mathbf{X}_n$ may not be Gaussian, its log-likelihood still performs well as a loss function in both theoretical and practical aspects. Sections 2.1 and 2.2 study the properties of the condition number constrained estimators of $\boldsymbol{\Theta}_0$, in the absence and presence of an L_1 penalty, respectively. Proofs of the results are given in the supplementary material.

2.1. Condition number constrained estimator: without L_1 penalty

Let $\boldsymbol{\Sigma}_0 = W_0 \Gamma_0 W_0$, where Γ_0 is the true correlation matrix and W_0 is the diagonal matrix of the true standard deviations. If $\Omega_0 = \Gamma_0^{-1}$, then $\boldsymbol{\Theta}_0 = W_0^{-1} \Omega_0 W_0^{-1}$. Let $\mathbf{S}_n = \widehat{W} \mathbf{R}_n \widehat{W}$, where \widehat{W}^2 is the diagonal matrix with the same diagonal as \mathbf{S}_n and \mathbf{R}_n is the sample correlation matrix.

We propose an estimator of the precision matrix

$$\widehat{\boldsymbol{\Theta}}_{\text{prop-1}} = \widehat{W}^{-1} \widetilde{\Omega}_{\kappa_n} \widehat{W}^{-1}, \quad (2.1)$$

where $\widetilde{\Omega}_{\kappa_n}$ is the solution to

$$\begin{cases} \underset{\Omega > 0}{\text{minimize}} & -\log\{\det(\Omega)\} + \text{tr}(\mathbf{R}_n \Omega), \\ \text{subject to} & \text{cond}(\Omega) \leq \kappa_n, \end{cases} \quad (2.2)$$

with a tuning parameter $\kappa_n \geq 1$.

Won et al. (2013) developed a well-conditioned estimator $\widehat{\boldsymbol{\Sigma}}_{\text{WLKR}}$ of $\boldsymbol{\Sigma}_0$. The inverse of their estimator, $\widehat{\boldsymbol{\Theta}}_{\text{WLKR}} = \widehat{\boldsymbol{\Sigma}}_{\text{WLKR}}^{-1}$, is also well-conditioned and can

be obtained by solving

$$\begin{cases} \underset{\Theta > 0}{\text{minimize}} & -\log\{\det(\Theta)\} + \text{tr}(\mathbf{S}_n \Theta), \\ \text{subject to} & \text{cond}(\Theta) \leq \kappa_n, \end{cases} \quad (2.3)$$

where $\kappa_n \geq 1$ is a tuning parameter. The difference between $\hat{\Theta}_{\text{prop-1}}$ in (2.1) and $\hat{\Theta}_{\text{WLKR}}$ in (2.3) is that $\hat{\Theta}_{\text{prop-1}}$ is a correlation-based estimator while $\hat{\Theta}_{\text{WLKR}}$ regularizes the precision matrix directly.

From Won et al. (2013), $\tilde{\Omega}_{\kappa_n}^{-1}$ in (2.2) truncates the eigenvalues of \mathbf{R}_n . However, because of the condition number constraint, $\tilde{\Omega}_{\kappa_n}^{-1}$ and \mathbf{R}_n have different asymptotic behaviors. We examine the asymptotic properties of $\hat{\Theta}_{\text{prop-1}}$. The following conditions will be involved.

- A1. $\lim_{n \rightarrow \infty} \log(p_n)/n = 0$, $\max_{1 \leq j \leq p_n} E[e^{t\{X_{1,j} - E(X_{1,j})\}^2}] < C$ for $|t| < c$, Σ_0 is diagonal, and $\kappa_n = 1$, where $C \in (0, \infty)$ and $c \in (0, \infty)$.
- A2. $\lim_{n \rightarrow \infty} p_n^{4/\beta}/n = 0$, $\max_{1 \leq j \leq p_n} E\{|X_{1,j} - E(X_{1,j})|^\beta\} < C$, Σ_0 is diagonal, and $\kappa_n = 1$, where $\beta \in [4, \infty)$ and $C \in (0, \infty)$.
- A3. $\lim_{n \rightarrow \infty} p_n/n = 0$, $\max_{1 \leq j \leq p_n} E\{|X_{1,j} - E(X_{1,j})|^4\} < C$, $\|\Sigma_0^{-1/2}\|_\infty < C$ with constant $C \in (0, \infty)$, and $\liminf_{n \rightarrow \infty} \{\kappa_n - \text{cond}(\Omega_0)\} > 0$. For any $i = 1, \dots, n$, $\{\mathbf{e}_{j,p_n}^T \Sigma_0^{-1/2}(\mathbf{X}_i - \boldsymbol{\mu}_0) : j = 1, \dots, p_n\}$ are i.i.d. random variables.

Theorem 1 (consistency of $\hat{\Theta}_{\text{prop-1}}$). *Suppose $\mathbf{X}_1, \dots, \mathbf{X}_n \in \mathbb{R}^{p_n}$ are i.i.d. with mean vector $\boldsymbol{\mu}_0$ and covariance matrix Σ_0 , and $0 < c \leq \lambda_{\min}(\Sigma_0) \leq \lambda_{\max}(\Sigma_0) \leq C < \infty$. For $\hat{\Theta}_{\text{prop-1}}$ in (2.1), under Condition A1 or A2 or A3, we have $\|\hat{\Theta}_{\text{prop-1}} - \Theta_0\|_2 \xrightarrow{P} 0$ and $\|\hat{\Theta}_{\text{prop-1}}^{-1} - \Sigma_0\|_2 \xrightarrow{P} 0$ as $n \rightarrow \infty$.*

The results of Theorem 1 also hold for $\hat{\Theta}_{\text{WLKR}}$ if we replace Conditions A1, A2, and A3 respectively by the following.

- A1*. $\lim_{n \rightarrow \infty} \log(p_n)/n = 0$, $\max_{1 \leq j \leq p_n} E[e^{t\{X_{1,j} - E(X_{1,j})\}^2}] < C$ for $|t| < c$, $\Sigma_0 = \mathbf{I}_{p_n}$, and $\kappa_n = 1$, where $C \in (0, \infty)$ and $c \in (0, \infty)$.
- A2*. $\lim_{n \rightarrow \infty} p_n^{4/\beta}/n = 0$, $\max_{1 \leq j \leq p_n} E\{|X_{1,j} - E(X_{1,j})|^\beta\} < C$, $\Sigma_0 = \mathbf{I}_{p_n}$, and $\kappa_n = 1$, where $\beta \in [4, \infty)$ and $C \in (0, \infty)$.
- A3*. $\lim_{n \rightarrow \infty} p_n/n = 0$, $\max_{1 \leq j \leq p_n} E\{|X_{1,j} - E(X_{1,j})|^4\} < C$, $\|\Sigma_0^{-1/2}\|_\infty < C$ with constant $C \in (0, \infty)$, and $\liminf_{n \rightarrow \infty} \{\kappa_n - \text{cond}(\Theta_0)\} > 0$. For any $i = 1, \dots, n$, $\{\mathbf{e}_{j,p_n}^T \Sigma_0^{-1/2}(\mathbf{X}_i - \boldsymbol{\mu}_0) : j = 1, \dots, p_n\}$ are i.i.d. random variables.

Condition A3 has been considered by Bai and Yin (1993) and El Karoui (2009), and holds when $\mathbf{X}_1, \dots, \mathbf{X}_n$ are, for example, normal. Comparing Conditions A1

and A2 with A1* and A2*, in high-dimensional settings, $\widehat{\Theta}_{\text{prop-1}}$ is consistent for Θ_0 when Σ_0 is diagonal, while the consistency of $\widehat{\Theta}_{\text{WLKR}}$ requires $\Sigma_0 = \mathbf{I}_{p_n}$. Conditions A1–A3 are restrictive. Theorem 2 details situations in which $\widehat{\Theta}_{\text{prop-1}}$ is inconsistent.

Denote by \mathbb{F}^{Γ_0} the empirical spectral distribution of Γ_0 . Suppose \mathbb{F}^{Γ_0} converges to a probability distribution function F_0 weakly as $n \rightarrow \infty$, and let

$$\begin{aligned} l_{\min} &= \inf\{x : F_0(x) > 0\}, & l_{\max} &= \sup\{x : F_0(x) < 1\}, \\ c_{\min} &= \inf\{F_0(x) : F_0(x) > 0\}, & c_{\max} &= \sup\{F_0(x) : F_0(x) < 1\}. \end{aligned} \quad (2.4)$$

Denote by $\mathbb{F}^{\mathbf{R}_n}$ the empirical spectral distribution of \mathbf{R}_n . From Theorem 1 of El Karoui (2009), if $\mathbf{X}_1, \dots, \mathbf{X}_n$ are i.i.d. random vectors, $\{e_{j,p_n}^T \Sigma_0^{-1/2} (\mathbf{X}_i - \boldsymbol{\mu}_0) : j = 1, \dots, p_n\}$ are i.i.d. random variables for any $i = 1, \dots, n$, \mathbb{F}^{Γ_0} converges to F_0 weakly, $\max_{1 \leq j \leq p_n} E\{|X_{1,j} - E(X_{1,j})|^\beta\} < C$, $\|\Sigma_0^{-1/2}\|_\infty < C$, $\|\Gamma_0\|_2 < C$, and $\lim_{n \rightarrow \infty} p_n/n = y$ with constants $\beta \in (4, \infty)$, $C \in (0, \infty)$, and $y \in (0, \infty)$, then $\mathbb{F}^{\mathbf{R}_n}$ converges weakly to a distribution function F in probability. We provide regularity conditions that are needed in Theorem 2.

- B1. $\|\widehat{W}^2 - W_0^2\|_2 = o_P(1)$, $\lim_{n \rightarrow \infty} p_n/n = \infty$, and $F_0 \neq \mathbf{I}_{[C, \infty)}$ for any $C \in [0, \infty)$.
- B2. $\|\widehat{W}^2 - W_0^2\|_2 = o_P(1)$ and $|\min\{\kappa_n, \text{cond}(\mathbf{R}_n)\} - \text{cond}(\Gamma_0)| \rightarrow 0$ in probability as $n \rightarrow \infty$.
- B3. $\{e_{j,p_n}^T \Sigma_0^{-1/2} (\mathbf{X}_i - \boldsymbol{\mu}_0) : j = 1, \dots, p_n\}$ are i.i.d. random variables, for any $i = 1, \dots, n$. $0 < l_{\min} < l_{\max} < \infty$, $\max_{1 \leq j \leq p_n} E\{|X_{1,j} - E(X_{1,j})|^\beta\} < C$, $\|\Sigma_0^{-1/2}\|_\infty < C$, $\lim_{n \rightarrow \infty} p_n/n = y$, $\lim_{n \rightarrow \infty} \kappa_n = l_{\max}/l_{\min}$, and $F_0 \neq F \mathbf{I}_{[l_{\min}, l_{\max})} + \mathbf{I}_{[l_{\max}, \infty)}$, where l_{\min} and l_{\max} are defined in (2.4), $\beta \in (4, \infty)$, $C \in (0, \infty)$, and $y \in (0, \infty)$ are constants and F is the limit that $\mathbb{F}^{\mathbf{R}_n}$ converges weakly to in probability.

Theorem 2 (inconsistency of $\widehat{\Theta}_{\text{prop-1}}$). *Suppose $\mathbf{X}_1, \dots, \mathbf{X}_n \in \mathbb{R}^{p_n}$ are i.i.d. with mean vector $\boldsymbol{\mu}_0$ and covariance matrix Σ_0 , $0 < c \leq \lambda_{\min}(\Sigma_0) \leq \lambda_{\max}(\Sigma_0) \leq C < \infty$. If \mathbb{F}^{Γ_0} converges to a probability distribution function F_0 weakly, for $\widehat{\Theta}_{\text{prop-1}}$ in (2.1), under Condition B1 or B2 or B3, we have $\|\widehat{\Theta}_{\text{prop-1}} - \Theta_0\|_2 \rightarrow 0$ and $\|\widehat{\Theta}_{\text{prop-1}}^{-1} - \Sigma_0\|_2 \rightarrow 0$ in probability.*

From Condition B1, $F_0 \neq \mathbf{I}_{[C, \infty)}$ for any $C \in [0, \infty)$ excludes the case of Σ_0 being diagonal. In this situation, if p_n is much larger than n , then $\widehat{\Theta}_{\text{prop-1}}$ is not a consistent estimator. Condition B2 implies that a well-selected tuning parameter κ_n is very important for the consistency of $\widehat{\Theta}_{\text{prop-1}}$. For Condition

B3, $F_0 \neq F I_{[l_{\min}, l_{\max})} + I_{[l_{\max}, \infty)}$ is satisfied in many situations. For example, from Silverstein and Choi (1995), if $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{\text{i.i.d.}}{\sim} N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$, \mathbb{F}^{Γ_0} converges to F_0 weakly, $\lim_{n \rightarrow \infty} p_n/n = y \in (0, 1)$ and some other regularity conditions hold, then F has a continuous density function on $(0, \infty)$. If F_0 does not have a continuous density, for example, F_0 is discrete, and $c_{\min} + c_{\max} < 1$ with c_{\min} and c_{\max} defined in (2.4), then condition $F_0 \neq F I_{[l_{\min}, l_{\max})} + I_{[l_{\max}, \infty)}$ holds. In this situation, $\hat{\boldsymbol{\Theta}}_{\text{prop-1}}$ is not consistent for $\boldsymbol{\Theta}_0$ in high-dimensional cases.

The results of Theorem 2 also hold for $\hat{\boldsymbol{\Theta}}_{\text{WLKR}}$ if we replace $\Omega_0, \Gamma_0, \mathbf{R}_n$, and $\hat{\boldsymbol{\Theta}}_{\text{prop-1}}$ with $\boldsymbol{\Theta}_0, \boldsymbol{\Sigma}_0, \mathbf{S}_n$, and $\hat{\boldsymbol{\Theta}}_{\text{WLKR}}$, respectively, in the conditions of Theorem 2.

2.2. Condition number constrained estimator: with L_1 penalty

To obtain consistent estimators of the precision matrix in high-dimensional cases, several estimation methods have been developed (see, e.g., Lam and Fan (2009) and Rothman et al. (2008)). The commonly used approach is the penalized log-likelihood method. In this section, we develop an estimator of $\boldsymbol{\Theta}_0$ with the condition number constraint and L_1 penalty, and study its asymptotic property.

By adding the L_1 penalty of Ω to the objective function in (2.2), we propose the estimator,

$$\hat{\boldsymbol{\Theta}}_{\text{prop-2}} = \widehat{W}^{-1} \widehat{\Omega}_{\mu_n, \kappa_n} \widehat{W}^{-1}, \tag{2.5}$$

where $\widehat{\Omega}_{\mu_n, \kappa_n}$ solves the optimization problem,

$$\begin{cases} \underset{\Omega \succ 0}{\text{minimize}} & -\log\{\det(\Omega)\} + \text{tr}(\mathbf{R}_n \Omega) + \mu_n |\Omega|_1, \\ \text{subject to} & \text{cond}(\Omega) \leq \kappa_n, \end{cases} \tag{2.6}$$

with $\mu_n > 0$ and $\kappa_n \geq 1$ the tuning parameters and $|\cdot|_1$ the matrix off-diagonal elementwise L_1 norm defined in (1.1). From the proof of Lemma 3 in Ravikumar et al. (2011), since the objective function in (2.6) is strictly convex for $\Omega \succ 0$ (Ravikumar et al. (2011)) and $\{\Omega \succ 0, \text{cond}(\Omega) \leq \kappa_n\}$ is a convex constraint, there exists a unique solution to (2.6).

A similar estimator without the condition number constraint developed in Rothman et al. (2008) is given by

$$\hat{\boldsymbol{\Theta}}_{\text{RBLZ}} = \widehat{W}^{-1} \widehat{\Omega}_{\text{RBLZ}} \widehat{W}^{-1}, \tag{2.7}$$

where

$$\widehat{\Omega}_{\text{RBLZ}} = \arg \min_{\Omega \succ 0} [-\log\{\det(\Omega)\} + \text{tr}(\mathbf{R}_n \Omega) + \mu_n |\Omega|_1],$$

with a tuning parameter $\mu_n > 0$. The convergence rate of $\widehat{\Theta}_{\text{RBLZ}}$ has been demonstrated in Rothman et al. (2008) under normal $\mathbf{X}_1, \dots, \mathbf{X}_n$. In Theorem 3, we re-examine the consistency of $\widehat{\Theta}_{\text{RBLZ}}$ under the exponential tail assumption. The convergence rate of $\widehat{\Theta}_{\text{RBLZ}}$ is also established under the polynomial tail condition. In the following, for sets $T, T' \subseteq \{1, \dots, p_n\} \times \{1, \dots, p_n\}$, let $(\Gamma_0 \otimes \Gamma_0)_{TT'}$ denote the $|T| \times |T'|$ submatrix of $\Gamma_0 \otimes \Gamma_0$ with rows and columns indexed by T and T' , respectively (see Section 3.1 in Ravikumar et al. (2011)). Specifically, if $T = \{(i_u, j_u) : u = 1, \dots, h\}$ and $T' = \{(i'_v, j'_v) : v = 1, \dots, h'\}$, then $e_{u,h}^T (\Gamma_0 \otimes \Gamma_0)_{TT'} e_{v,h'} = \Gamma_0(i_u, i'_v) \Gamma_0(j_u, j'_v)$ for $u = 1, \dots, h$ and $v = 1, \dots, h'$.

For the next result, let $s_n = |\{(i, j) : i \neq j \text{ and } \Theta_0(i, j) \neq 0\}|$ and $t_n = \max_{i=1, \dots, p_n} |\{j = 1, \dots, p_n : \Theta_0(i, j) \neq 0\}|$ and consider the following conditions

- C1. (exponential tail condition) $\max_{1 \leq j \leq p_n} E[e^{t\{X_{1,j} - E(X_{1,j})\}^2}] < C$ for $|t| < c$ with certain constants $C \in (0, \infty)$ and $c \in (0, \infty)$, $\mu_n \asymp \{\log(p_n)/n\}^{1/2}$ and $r_n = o(1)$ where $r_n = \min(1 + s_n, t_n^2) \log(p_n)/n$.
- C2. (polynomial tail condition) $\max_{1 \leq j \leq p_n} E\{|X_{1,j} - E(X_{1,j})|^\beta\} < C$ for certain constants $\beta \in [4, \infty)$ and $C \in (0, \infty)$, $\mu_n \asymp p_n^{2\tau/\beta}/n^{1/2}$ and $r_n = o(1)$ where $r_n = \min(1 + s_n, t_n^2) p_n^{4\tau/\beta}/n$ and $\tau \in (2, \infty)$ is a constant.

Theorem 3. *Suppose $\mathbf{X}_1, \dots, \mathbf{X}_n \in \mathbb{R}^{p_n}$ are i.i.d. random vectors with covariance matrix Σ_0 such that $\lambda_{\min}(\Sigma_0) \geq c > 0$ and $\|\Sigma_0\|_\infty \leq C < \infty$. Let $S = \{(i, j) : \Theta_0(i, j) \neq 0\}$ and assume $\|(\Gamma_0 \otimes \Gamma_0)_{SS}\|_\infty^{-1} \leq C < \infty$ and $\max_{e \in \{(i,j) : \Theta_0(i,j)=0\}} \|(\Gamma_0 \otimes \Gamma_0)_{eS}\|_1 \leq c$ with constant $c \in (0, 1)$. For $\widehat{\Theta}_{\text{RBLZ}}$ in (2.7), if either C1 or C2 holds, we have $\|\widehat{\Theta}_{\text{RBLZ}} - \Theta_0\|_2^2 = O_P(r_n) = \|\widehat{\Theta}_{\text{RBLZ}}^{-1} - \Sigma_0\|_2^2$.*

Rothman et al. (2008) demonstrated that the convergence rate of $\widehat{\Theta}_{\text{RBLZ}}$ under the L_2 norm is $\{(1 + s_n) \log(p_n)/n\}^{1/2}$ by assuming that $\mathbf{X}_1, \dots, \mathbf{X}_n$ are multivariate normal and the eigenvalues of Σ_0 are bounded away from 0 and ∞ . If, in Θ_0 , the maximum number of non-zeros per row is large relative to the total number of non-zero off-diagonal elements, $s_n = O(t_n^2)$, then the convergence rate of $\widehat{\Theta}_{\text{RBLZ}}$ in Theorem 3 under the exponential tail condition is equivalent to that in Rothman et al. (2008). However, when $t_n^2 = o(s_n)$, our result provides a faster convergence rate while requiring stronger conditions than that in Rothman et al. (2008). Particularly, conditions $\|(\Gamma_0 \otimes \Gamma_0)_{SS}\|_\infty^{-1} \leq C < \infty$ and $\max_{e \in \{(i,j) : \Theta_0(i,j)=0\}} \|(\Gamma_0 \otimes \Gamma_0)_{eS}\|_1 \leq c$ in Theorem 3 are adopted from Ravikumar et al. (2011), the latter of which is regarded as the mutual incoherence or irrepresentability condition.

We now study the asymptotic properties of $\widehat{\Theta}_{\text{prop-2}}$.

Theorem 4 (consistency of $\widehat{\Theta}_{\text{prop-2}}$). *For $\widehat{\Theta}_{\text{prop-2}}$ in (2.5), under the conditions in Theorem 3 and $\liminf_{n \rightarrow \infty} \{\kappa_n - \text{cond}(\Omega_0)\} > 0$,*

$$\begin{aligned} \lim_{n \rightarrow \infty} P(\widehat{\Theta}_{\text{prop-2}} = \widehat{\Theta}_{\text{RBLZ}}) &= 1, \\ \|\widehat{\Theta}_{\text{prop-2}} - \Theta_0\|_2^2 &= O_P(r_n) = \|\widehat{\Theta}_{\text{prop-2}}^{-1} - \Sigma_0\|_2^2. \end{aligned}$$

Comparing the results of Theorems 1, 2, and 4, we can see the effect of the L_1 penalty. Under the condition that $\max(r_n, p_n/n) = o(1)$, $\widehat{\Theta}_{\text{prop-1}}$ and $\widehat{\Theta}_{\text{prop-2}}$ are consistent. In high-dimensional cases with $\lim_{n \rightarrow \infty} p_n/n > 0$ and $r_n = o(1)$, if Σ_0 is not diagonal, then, under some regularity conditions, $\widehat{\Theta}_{\text{prop-2}}$ converges to Θ_0 but $\widehat{\Theta}_{\text{prop-1}}$ may not. The L_1 penalty of Ω in the objective function of (2.6) is necessary for the convergence of $\widehat{\Theta}_{\text{prop-2}}$ in high-dimensional settings.

Although $\widehat{\Theta}_{\text{prop-2}}$ is asymptotically equivalent to $\widehat{\Theta}_{\text{RBLZ}}$, we find in numerical studies in Sections 5 and 6 that our proposed method performs better in finite-sample situations while effectively controlling the condition number of the estimate.

Another competitive estimator is the graphical lasso estimator (Friedman, Hastie and Tibshirani (2008)) defined as

$$\widehat{\Theta}_{\text{GLasso}} = \arg \min_{\Theta \succ 0} [-\log\{\det(\Theta)\} + \text{tr}(\mathbf{S}_n \Theta) + \mu_n |\Theta|_1], \tag{2.8}$$

with a tuning parameter $\mu_n > 0$. Consistency of $\widehat{\Theta}_{\text{GLasso}}$ has been demonstrated by Rothman et al. (2008) and Ravikumar et al. (2011) under different conditions. The convergence rate of $\widehat{\Theta}_{\text{GLasso}}$ under the L_2 norm derived in Rothman et al. (2008) is $\{(p_n + s_n) \log(p_n)/n\}^{1/2}$ under normality. Comparing the rate of $\widehat{\Theta}_{\text{GLasso}}$ in Rothman et al. (2008) with that of $\widehat{\Theta}_{\text{prop-2}}$ in Theorem 4 under Condition C1, we find that the two rates are equivalent if $p_n = O(s_n)$ and $s_n = O(t_n^2)$, while $\widehat{\Theta}_{\text{prop-2}}$ converges faster if $t_n^2 = o(p_n + s_n)$ or $s_n = o(p_n)$. In Ravikumar et al. (2011), the derived convergence rate of $\widehat{\Theta}_{\text{GLasso}}$ is $\{\min(p_n + s_n, t_n^2) \log(p_n)/n\}^{1/2}$ under the exponential tail assumption and $\{\min(p_n + s_n, t_n^2) p_n^{4\tau/\beta}/n\}^{1/2}$ under the polynomial tail assumption. Under either the exponential tail or polynomial tail assumption, the convergence rate of $\widehat{\Theta}_{\text{prop-2}}$ is equivalent to that of $\widehat{\Theta}_{\text{GLasso}}$ in Ravikumar et al. (2011) if $p_n = O(s_n)$ or $t_n^2 = O(s_n)$, while the rate of $\widehat{\Theta}_{\text{prop-2}}$ is sharper if $s_n = o(t_n^2)$ and $s_n = o(p_n)$. To sum up, compared with the convergence rates of $\widehat{\Theta}_{\text{GLasso}}$ obtained in past work (Rothman et al. (2008); Ravikumar et al. (2011)), the rate of convergence for $\widehat{\Theta}_{\text{prop-2}}$ is faster under certain situations, for example, in cases where Θ_0 is sparse and the total number of non-zero off-

diagonal elements in Θ_0 is small relative to the maximum number of non-zeros per row. Under some other situations, where s_n dominates p_n or t_n^2 , $\hat{\Theta}_{\text{prop-2}}$ and $\hat{\Theta}_{\text{GLasso}}$ converge to Θ_0 at the same rate.

3. Algorithm for Solving (2.6)

Numerically solving the optimization problem (2.6) is a non-trivial task. Different algorithms for penalized sparse precision matrix estimation have been developed (Boyd et al. (2010); Friedman, Hastie and Tibshirani (2008)). For example, Boyd et al. (2010) proposed an ADMM (alternating direction method of multipliers) algorithm, which can be used to solve the optimization problems for $\hat{\Theta}_{\text{RBLZ}}$ and $\hat{\Theta}_{\text{GLasso}}$. Based on the ADMM algorithm, we develop an algorithm for our estimator. The problem (2.6) is equivalent to the optimization problem

$$\begin{cases} \text{minimize}_{\Omega \succ 0, Z} & -\log\{\det(\Omega)\} + \text{tr}(\mathbf{R}_n \Omega) + \mu_n |Z|_1, \\ \text{subject to} & \text{cond}(\Omega) \leq \kappa_n, \Omega = Z. \end{cases} \quad (3.1)$$

To deal with (3.1), we minimize the corresponding scaled augmented Lagrangian,

$$\begin{cases} \text{minimize}_{\Omega \succ 0, Z} & L_\rho(\Omega, Z; U), \\ \text{subject to} & \text{cond}(\Omega) \leq \kappa_n, \Omega = Z, \end{cases} \quad (3.2)$$

where

$$\begin{aligned} L_\rho(\Omega, Z; U) = & -\log\{\det(\Omega)\} + \text{tr}(\mathbf{R}_n \Omega) + \mu_n |Z|_1 \\ & + \frac{\rho}{2} \|\Omega - Z + U\|_F^2 - \frac{\rho}{2} \|U\|_F^2, \end{aligned} \quad (3.3)$$

is the scaled augmented Lagrange function, $\rho \in (0, \infty)$ is an arbitrary constant, and U is the Lagrange multiplier. The objective functions and constraints in both (3.1) and (3.2) are convex, and therefore there exist unique solutions to the two optimization problems (see, e.g., the proof of Lemma 3 in Ravikumar et al. (2011)). The problems (2.6), (3.1), and (3.2) are equivalent, so we use (3.2).

Motivated by the ADMM algorithm (Boyd et al. (2010)), we calculate the limit of $Z^{(i)}$ as the solution to (3.2) by iterations (with $i = 1, 2, \dots$) of three steps until convergence

$$\begin{aligned} \text{Step 1 : } \Omega^{(i)} & \leftarrow \arg \min_{\Omega \succ 0, \text{cond}(\Omega) \leq \kappa_n} L_\rho(\Omega, Z^{(i-1)}; U^{(i-1)}), \\ \text{Step 2 : } Z^{(i)} & \leftarrow \arg \min_Z L_\rho(\Omega^{(i)}, Z; U^{(i-1)}), \\ \text{Step 3 : } U^{(i)} & \leftarrow U^{(i-1)} + \Omega^{(i)} - Z^{(i)}. \end{aligned}$$

The criterion for declaring algorithmic convergence is

$$\frac{\sum_{j=1}^{p_n} \sum_{k=1}^{p_n} |\Omega^{(i+1)}(j, k) - \Omega^{(i)}(j, k)|}{\sum_{j=1}^{p_n} \sum_{k=1}^{p_n} |\Omega^{(i)}(j, k)|} \leq 10^{-4}.$$

Issues on the global convergence of the ADMM algorithm can be found in Boyd et al. (2010) (see Section 3.2 therein for details). In practice, we use the zero matrix as the initial values for $Z^{(0)}$ and $U^{(0)}$, and set $\rho = 1$ for each iteration to control the step size.

From Boyd et al. (2010), the optimization problem for $\hat{\Theta}_{\text{RBLZ}}$ can be solved using the ADMM algorithm by iterations of three steps with an explicit solution for each step. The ADMM algorithm for $\hat{\Theta}_{\text{prop-2}}$ here differs from that for $\hat{\Theta}_{\text{RBLZ}}$ in the sense that we consider the constraint $\{\Omega \succ 0, \text{cond}(\Omega) \leq \kappa_n\}$ in Step 1 while they use $\Omega \succ 0$.

Next, we calculate the solutions for Steps 1 and 2 in our algorithm. For Step 2, from (3.3),

$$\begin{aligned} Z^{(i)} &= \arg \min_Z L_\rho(\Omega^{(i)}, Z; U^{(i-1)}) \\ &= \arg \min_Z \left\{ \mu_n |Z|_1 + \frac{\rho}{2} \|\Omega^{(i)} - Z + U^{(i-1)}\|_F^2 \right\} \\ &= \arg \min_Z \left[\frac{1}{2} \|Z - \{\Omega^{(i)} + U^{(i-1)}\}\|_F^2 + \frac{\mu_n}{\rho} |Z|_1 \right]. \end{aligned} \tag{3.4}$$

The last optimization problem in (3.4) is similar to problem (1) in Xue, Ma and Zou (2012), which has a closed-form solution by soft-thresholding (see Paragraph 2 of Section 1 in Xue, Ma and Zou (2012) for details). By arguments similar to those in Xue, Ma and Zou (2012), we can show that there also exists a closed-form solution to (3.4) based on soft-thresholding: for $j, k = 1, \dots, p_n$,

$$Z^{(i)}(j, k) = \begin{cases} A^{(i)}(j, k), & \text{if } j = k, \\ \text{sign}\{A^{(i)}(j, k)\} \max\{|A^{(i)}(j, k)| - \mu_n/\rho, 0\}, & \text{otherwise,} \end{cases}$$

where $A^{(i)} = \Omega^{(i)} + U^{(i-1)}$. For Step 1, the solution is not straightforward due to the constraint $\{\Omega \succ 0, \text{cond}(\Omega) \leq \kappa_n\}$. To obtain $\Omega^{(i)}$, we now propose a method, the proof of which is available in the online supplementary material.

Proposition 1. *For the optimization problem in Step 1, let VDV^T be the eigen-decomposition of $\mathbf{R}_n/\rho - Z^{(i-1)} + U^{(i-1)}$ with $D = \text{diag}(d_1, \dots, d_{p_n})$ and $d_1 \geq \dots \geq d_{p_n}$. Let $\delta_j = -d_j/2 + \sqrt{d_j^2/4 + 1/\rho}$ for $j = 1, \dots, p_n$,*

$$\tilde{D} = \begin{cases} \text{diag}(\delta_1, \dots, \delta_{p_n}), & \text{if } \delta_{p_n}/\delta_1 \leq \kappa_n, \\ \text{diag}(\tilde{d}_1, \dots, \tilde{d}_{p_n}), & \text{if } \delta_{p_n}/\delta_1 > \kappa_n, \end{cases}$$

where $\tilde{d}_j = \min\{\max(\tau_0, \delta_j), \kappa_n \tau_0\}$ and

$$\tau_0 = \frac{\left[-\rho \left(\sum_{j=1}^{\alpha_0} d_j + \kappa_n \sum_{j=\beta_0}^{p_n} d_j \right) + \left\{ \rho^2 \left(\sum_{j=1}^{\alpha_0} d_j + \kappa_n \sum_{j=\beta_0}^{p_n} d_j \right)^2 + 4\rho(\alpha_0 + \kappa_n^2 p_n - \kappa_n^2 \beta_0 + \kappa_n^2)(\alpha_0 + p_n - \beta_0 + 1) \right\}^{1/2} \right]}{2\rho(\alpha_0 + \kappa_n^2 p_n - \kappa_n^2 \beta_0 + \kappa_n^2)},$$

with α_0 the largest index in $\{1, \dots, p_n\}$ such that $\tau_0 > \delta_{\alpha_0}$ and β_0 the smallest index in $\{1, \dots, p_n\}$ such that $\kappa_n \tau_0 < \delta_{\beta_0}$. Then, the solution to Step 1 is $\Omega^{(i)} = V\tilde{D}V^T$. The quantities α_0 and β_0 can be found in $O(p_n)$ operations.

Since Step 1 also requires the eigendecomposition of $\mathbf{R}_n/\rho - Z^{(i-1)} + U^{(i-1)}$ which takes $O(p_n^3)$ operations, the number of operations for each iteration of the algorithm is $O(p_n^3)$. To calculate $\hat{\Theta}_{\text{RBLZ}}$ or $\hat{\Theta}_{\text{GLasso}}$, the ADMM algorithm also needs $O(p_n^3)$ operations.

4. Tuning Parameter Selection

This section illustrates a data-driven method to select the tuning parameters for $\hat{\Theta}_{\text{prop-2}}$. In practice, we choose κ_n and μ_n in an iterative way. At each step, one of them is fixed and the other one is updated, with cross validation used for choosing κ_n and BIC for μ_n . Specifically, for a fixed μ_n , we divide the data into k folds and choose κ_n by minimizing

$$\text{CV}(\kappa_n; \mu_n) = \sum_{i=1}^k \frac{n}{2k} [-\log\{\det(\hat{\Omega}_{\mu_n, \kappa_n}^{[-i]})\} + \text{tr}(\mathbf{R}_n^{[i]} \hat{\Omega}_{\mu_n, \kappa_n}^{[-i]})],$$

where $\mathbf{R}_n^{[i]}$ is the sample correlation matrix based on the i th fold and $\hat{\Omega}_{\mu_n, \kappa_n}^{[-i]}$ is the estimate of Ω_0 calculated with all observations except those in the i th fold. Given κ_n , we choose μ_n that minimizes the BIC function

$$\begin{aligned} \text{BIC}(\mu_n; \kappa_n) &= -n \log\{\det(\hat{\Omega}_{\mu_n, \kappa_n})\} + n \text{tr}(\mathbf{R}_n \hat{\Omega}_{\mu_n, \kappa_n}) \\ &\quad + \log(n) \sum_{1 \leq i \leq j \leq p_n} \mathbf{I}\{\hat{\Omega}_{\mu_n, \kappa_n}(i, j) \neq 0\}. \end{aligned}$$

The details for selecting κ_n and μ_n are as follows.

Step I : Initialize κ_n^0 .

Step II : Repeat the following steps (with $i = 1, 2, \dots$) until convergence:

$$\begin{aligned}\mu_n^i &= \arg \min_{\mu_n} \text{BIC}(\mu_n; \kappa_n^{i-1}), \\ \kappa_n^i &= \arg \min_{\kappa_n} \text{CV}(\kappa_n; \mu_n^i),\end{aligned}$$

where the optimization problems are solved by grid search.

In the numerical studies in Sections 5 and 6, the initial value in Step I is $\kappa_n^0 = \infty$. We also observe that the averaged number of iterations for the numerical convergence of the algorithm in Step II is moderate in each situation, and does not increase as p_n increases.

5. Simulation Evaluation

Simulation studies were conducted to compare estimators $\widehat{\Theta}_{\text{prop-1}}$ in (2.1), $\widehat{\Theta}_{\text{prop-2}}$ in (2.5), $\widehat{\Theta}_{\text{WLKR}}$ in (2.3), $\widehat{\Theta}_{\text{RBLZ}}$ in (2.7), $\widehat{\Theta}_{\text{GLasso}}$ in (2.8), and $\widehat{\Theta}_{\text{banded}}$, with $n = 300$ and $p_n \in \{100, 200, 400\}$, where $\widehat{\Theta}_{\text{banded}}$ is the banded estimate of Θ_0 by Cholesky decomposition as defined in Bickel and Levina (2008) (see Section 2.2 therein). To generate data, we considered the following schemes

- I. $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{\text{i.i.d.}}{\sim} N(\mathbf{0}, \Sigma_0)$, where $\Sigma_0 = \text{diag}(10, \dots, 10, 0.01, \dots, 0.01)$. The proportion of the “high” eigenvalues is 80% among the p_n eigenvalues. Similar schemes were used in Won et al. (2013).
- II. $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{\text{i.i.d.}}{\sim} t_3(\mathbf{0}, \Sigma^*)$ with $(\Sigma^*)^{-1} = \text{diag}(3K, K, \dots, K)$ where K is a 50×50 matrix such that $K(i, j) = \text{I}(i = j) + 0.1 \text{I}(|i - j| = 1) + 0.4 \text{I}(|i - j| = 3)$ for $i, j = 1, \dots, 50$.
- III. $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{\text{i.i.d.}}{\sim} t_3(\mathbf{0}, \Sigma^*)$ with $\mathbf{e}_{i,p_n}^T (\Sigma^*)^{-1} \mathbf{e}_{j,p_n} = \{3 \text{I}(i = j) + 1.49 \text{I}(|i - j| = 1)\} \{1 - 1/2 \text{I}(i \leq p_n/2)\} \{1 - 1/2 \text{I}(j \leq p_n/2)\}$ for $i, j = 1, \dots, p_n$.

Here, $t_\nu(\mathbf{0}, \Sigma^*)$ denotes the multivariate t distribution with degree of freedom ν , location vector $\mathbf{0} \in \mathbb{R}^{p_n}$, and scale matrix Σ^* . Specifically, if random quantities $\mathbf{Y} \sim N(\mathbf{0}, \Sigma^*)$ and $w \sim \chi_\nu^2$ are independent, then $\mathbf{Y}/(w/\nu)^{1/2} \sim t_\nu(\mathbf{0}, \Sigma^*)$ (see, for example, Section 5.6 in DeGroot (2004)). Thus $\Sigma_0 = 3\Sigma^*$ in schemes II and III. For schemes I–III, the structures of Θ_0 are diagonal, block diagonal, and banded, respectively.

Monte Carlo simulations were replicated 400 times in each setting. For a generic estimator $\widehat{\Theta}$ of Θ_0 , we calculated the averaged losses $\|\widehat{\Theta} - \Theta_0\|_F$ and $\|\widehat{\Theta} - \Theta_0\|_2$. The selection performance was measured by the false positive rate (FPR) and false negative rate (FNR):

Table 1. Comparison of $\widehat{\Theta}_{\text{prop-1}}$, $\widehat{\Theta}_{\text{prop-2}}$, $\widehat{\Theta}_{\text{WLKR}}$, $\widehat{\Theta}_{\text{RBLZ}}$, $\widehat{\Theta}_{\text{GLasso}}$, and $\widehat{\Theta}_{\text{banded}}$ with data generated by scheme I. Each metric is averaged over 400 replications with the standard error shown in the bracket.

p_n	$\widehat{\Theta}$	$\ \widehat{\Theta} - \Theta_0\ _F$	$\ \widehat{\Theta} - \Theta_0\ _2$	FPR	FNR
100	$\widehat{\Theta}_{\text{prop-1}}$	36.96 (0.31)	18.81 (0.23)	0.00 (0.00)	0.00 (0.00)
	$\widehat{\Theta}_{\text{prop-2}}$	36.96 (0.31)	18.81 (0.23)	0.00 (0.00)	0.00 (0.00)
	$\widehat{\Theta}_{\text{WLKR}}$	140.22 (0.15)	31.46 (0.03)	1.00 (0.00)	0.00 (0.00)
	$\widehat{\Theta}_{\text{RBLZ}}$	36.97 (0.31)	18.82 (0.23)	0.00 (0.00)	0.00 (0.00)
	$\widehat{\Theta}_{\text{GLasso}}$	36.96 (0.31)	18.81 (0.23)	0.00 (0.00)	0.00 (0.00)
	$\widehat{\Theta}_{\text{banded}}$	49.92 (1.07)	24.66 (0.56)	0.02 (0.00)	0.00 (0.00)
200	$\widehat{\Theta}_{\text{prop-1}}$	52.64 (0.31)	21.42 (0.25)	0.00 (0.00)	0.00 (0.00)
	$\widehat{\Theta}_{\text{prop-2}}$	52.64 (0.31)	21.42 (0.25)	0.00 (0.00)	0.00 (0.00)
	$\widehat{\Theta}_{\text{WLKR}}$	631.40 (0.00)	99.83 (0.00)	1.00 (0.00)	0.00 (0.00)
	$\widehat{\Theta}_{\text{RBLZ}}$	52.65 (0.31)	21.42 (0.25)	0.00 (0.00)	0.00 (0.00)
	$\widehat{\Theta}_{\text{GLasso}}$	52.64 (0.31)	21.42 (0.25)	0.00 (0.00)	0.00 (0.00)
	$\widehat{\Theta}_{\text{banded}}$	72.62 (1.75)	28.26 (0.72)	0.01 (0.00)	0.00 (0.00)
400	$\widehat{\Theta}_{\text{prop-1}}$	74.40 (0.31)	23.48 (0.21)	0.00 (0.00)	0.00 (0.00)
	$\widehat{\Theta}_{\text{prop-2}}$	74.40 (0.31)	23.47 (0.21)	0.00 (0.00)	0.00 (0.00)
	$\widehat{\Theta}_{\text{WLKR}}$	892.90 (0.00)	99.83 (0.00)	1.00 (0.00)	0.00 (0.00)
	$\widehat{\Theta}_{\text{RBLZ}}$	74.41 (0.31)	23.48 (0.21)	0.00 (0.00)	0.00 (0.00)
	$\widehat{\Theta}_{\text{GLasso}}$	74.40 (0.31)	23.48 (0.21)	0.00 (0.00)	0.00 (0.00)
	$\widehat{\Theta}_{\text{banded}}$	113.92 (2.85)	33.39 (0.81)	0.01 (0.00)	0.00 (0.00)

$$\text{FPR} = \frac{|\{(i, j) : \Theta_0(i, j) = 0, \widehat{\Theta}(i, j) \neq 0\}|}{|\{(i, j) : \Theta_0(i, j) = 0\}|},$$

$$\text{FNR} = \frac{|\{(i, j) : \Theta_0(i, j) \neq 0, \widehat{\Theta}(i, j) = 0\}|}{|\{(i, j) : \Theta_0(i, j) \neq 0\}|}.$$

For $\widehat{\Theta}_{\text{prop-2}}$, the tuning parameters were selected as described in Section 4 with $k = 5$. To calculate $\widehat{\Theta}_{\text{prop-1}}$ and $\widehat{\Theta}_{\text{WLKR}}$, we solved (2.2) and (2.3) using the algorithm in Won et al. (2013) (see (8), Lemma 1, and Theorem 1). The tuning parameter κ_n of $\widehat{\Theta}_{\text{prop-1}}$ and $\widehat{\Theta}_{\text{WLKR}}$ were selected by the 5-fold CV as illustrated in Won et al. (2013). The graphical lasso algorithm in Friedman, Hastie and Tibshirani (2008) was applied to calculate $\widehat{\Theta}_{\text{RBLZ}}$ and $\widehat{\Theta}_{\text{GLasso}}$, where μ_n was selected by the 5-fold cross validation illustrated in Rothman et al. (2008). For these methods, κ_n was chosen from $\{1.4226^i : i = 0, 1, \dots, 29\}$ and μ_n from $\{0.02 \times 1.2390^i : i = 0, 1, \dots, 29\}$. The banding parameter for $\widehat{\Theta}_{\text{banded}}$ was selected by the random splitting method in Bickel and Levina (2008).

The averaged losses, FPR, and FNR of different precision matrix estimators are presented in Tables 1–3 corresponding to schemes I–III, respectively. For

Table 2. Comparison of $\hat{\Theta}_{\text{prop-1}}$, $\hat{\Theta}_{\text{prop-2}}$, $\hat{\Theta}_{\text{WLKR}}$, $\hat{\Theta}_{\text{RBLZ}}$, $\hat{\Theta}_{\text{GLasso}}$, and $\hat{\Theta}_{\text{banded}}$ with data generated by scheme II. Each metric is averaged over 400 replications with the standard error shown in the bracket.

p_n	$\hat{\Theta}$	$\ \hat{\Theta} - \Theta_0\ _F$	$\ \hat{\Theta} - \Theta_0\ _2$	FPR	FNR
100	$\hat{\Theta}_{\text{prop-1}}$	4.49 (0.02)	1.40 (0.01)	1.00 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{prop-2}}$	2.59 (0.03)	0.84 (0.01)	0.20 (0.00)	0.09 (0.00)
	$\hat{\Theta}_{\text{WLKR}}$	4.76 (0.02)	1.25 (0.01)	1.00 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{RBLZ}}$	2.93 (0.04)	0.94 (0.01)	0.13 (0.00)	0.10 (0.00)
	$\hat{\Theta}_{\text{GLasso}}$	2.94 (0.04)	0.89 (0.01)	0.15 (0.00)	0.11 (0.00)
	$\hat{\Theta}_{\text{banded}}$	5.06 (0.09)	1.62 (0.04)	0.07 (0.00)	0.15 (0.01)
200	$\hat{\Theta}_{\text{prop-1}}$	5.10 (0.02)	1.37 (0.01)	1.00 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{prop-2}}$	2.91 (0.03)	0.87 (0.01)	0.11 (0.00)	0.08 (0.00)
	$\hat{\Theta}_{\text{WLKR}}$	6.11 (0.03)	1.55 (0.01)	1.00 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{RBLZ}}$	3.31 (0.05)	0.97 (0.01)	0.08 (0.00)	0.10 (0.00)
	$\hat{\Theta}_{\text{GLasso}}$	3.84 (0.04)	1.07 (0.01)	0.09 (0.00)	0.11 (0.00)
	$\hat{\Theta}_{\text{banded}}$	5.76 (0.14)	1.77 (0.05)	0.05 (0.00)	0.10 (0.01)
400	$\hat{\Theta}_{\text{prop-1}}$	6.14 (0.02)	1.38 (0.01)	1.00 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{prop-2}}$	3.57 (0.03)	0.91 (0.01)	0.06 (0.00)	0.08 (0.00)
	$\hat{\Theta}_{\text{WLKR}}$	7.33 (0.03)	1.69 (0.00)	1.00 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{RBLZ}}$	3.91 (0.05)	0.97 (0.01)	0.05 (0.00)	0.09 (0.00)
	$\hat{\Theta}_{\text{GLasso}}$	4.65 (0.04)	1.20 (0.01)	0.06 (0.00)	0.11 (0.00)
	$\hat{\Theta}_{\text{banded}}$	7.17 (0.22)	1.94 (0.06)	0.03 (0.00)	0.06 (0.01)

scheme I, where Σ_0 is diagonal, $\hat{\Theta}_{\text{prop-1}}$, $\hat{\Theta}_{\text{prop-2}}$, $\hat{\Theta}_{\text{RBLZ}}$, and $\hat{\Theta}_{\text{GLasso}}$ perform comparably well and outperform $\hat{\Theta}_{\text{WLKR}}$ and $\hat{\Theta}_{\text{banded}}$. The results for $\hat{\Theta}_{\text{prop-1}}$, $\hat{\Theta}_{\text{prop-2}}$, $\hat{\Theta}_{\text{RBLZ}}$, and $\hat{\Theta}_{\text{GLasso}}$ are similar in Table 1, since the calculated estimates are all close to \widehat{W}^2 . From Table 2, by comparing the losses of the precision matrix estimators, we see that $\hat{\Theta}_{\text{prop-2}}$ outperforms the other estimators. From Table 3, $\hat{\Theta}_{\text{prop-2}}$ has smaller averaged losses than $\hat{\Theta}_{\text{prop-1}}$, $\hat{\Theta}_{\text{WLKR}}$, $\hat{\Theta}_{\text{GLasso}}$, and $\hat{\Theta}_{\text{banded}}$. Compared with $\hat{\Theta}_{\text{RBLZ}}$, when $p_n = 400$, $\hat{\Theta}_{\text{prop-2}}$ has a slightly larger loss under the L_2 norm. In the other cases, $\hat{\Theta}_{\text{prop-2}}$ performs better than $\hat{\Theta}_{\text{RBLZ}}$.

As suggested by one referee, the averaged condition numbers of $\hat{\Theta}_{\text{prop-2}}$ and $\hat{\Theta}_{\text{RBLZ}}$ were also compared under the same amount of L_1 regularization. After calculating $\hat{\Theta}_{\text{prop-2}}$ with the tuning parameters $\kappa_n = \hat{\kappa}_n$ and $\mu_n = \hat{\mu}_n$ selected by the data-driven method in Section 4, we calculated $\hat{\Theta}_{\text{RBLZ}}$ with μ_n equal to $\hat{\mu}_n$ instead of selected by CV, and denote the resulting estimator by $\hat{\Theta}_{\text{RBLZ}}^*$. Hence, $\hat{\Theta}_{\text{prop-2}}$ and $\hat{\Theta}_{\text{RBLZ}}^*$ have the same amount of L_1 regularization. However, for scheme I where $\Omega_0 = \mathbf{I}_{p_n}$, the data-driven choice of κ_n for $\hat{\Theta}_{\text{prop-2}}$ is exactly 1

Table 3. Comparison of $\hat{\Theta}_{\text{prop-1}}$, $\hat{\Theta}_{\text{prop-2}}$, $\hat{\Theta}_{\text{WLKR}}$, $\hat{\Theta}_{\text{RBLZ}}$, $\hat{\Theta}_{\text{GLasso}}$, and $\hat{\Theta}_{\text{banded}}$ with data generated by scheme III. Each metric is averaged over 400 replications with the standard error shown in the bracket.

p_n	$\hat{\Theta}$	$\ \hat{\Theta} - \Theta_0\ _F$	$\ \hat{\Theta} - \Theta_0\ _2$	FPR	FNR
100	$\hat{\Theta}_{\text{prop-1}}$	4.22 (0.03)	1.48 (0.02)	1.00 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{prop-2}}$	2.34 (0.04)	0.84 (0.01)	0.22 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{WLKR}}$	5.44 (0.03)	1.32 (0.01)	1.00 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{RBLZ}}$	3.02 (0.05)	0.95 (0.01)	0.14 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{GLasso}}$	3.63 (0.05)	1.01 (0.01)	0.18 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{banded}}$	4.28 (0.12)	1.46 (0.04)	0.04 (0.00)	0.09 (0.01)
200	$\hat{\Theta}_{\text{prop-1}}$	6.67 (0.05)	1.60 (0.01)	1.00 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{prop-2}}$	3.86 (0.06)	0.97 (0.01)	0.13 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{WLKR}}$	8.91 (0.04)	1.54 (0.01)	1.00 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{RBLZ}}$	4.58 (0.08)	1.03 (0.01)	0.10 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{GLasso}}$	5.93 (0.07)	1.14 (0.01)	0.12 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{banded}}$	6.90 (0.26)	1.91 (0.08)	0.04 (0.00)	0.06 (0.01)
400	$\hat{\Theta}_{\text{prop-1}}$	11.34 (0.05)	1.58 (0.01)	1.00 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{prop-2}}$	6.69 (0.06)	1.14 (0.01)	0.07 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{WLKR}}$	13.92 (0.04)	1.67 (0.00)	1.00 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{RBLZ}}$	7.00 (0.10)	1.11 (0.01)	0.06 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{GLasso}}$	9.43 (0.08)	1.25 (0.01)	0.07 (0.00)	0.00 (0.00)
	$\hat{\Theta}_{\text{banded}}$	9.81 (0.39)	2.08 (0.09)	0.02 (0.00)	0.03 (0.01)

in many replications, and hence $\hat{\Omega}_{\mu_n, \kappa_n} = \mathbf{I}_{p_n}$ for any $\mu_n > 0$. In this situation, $\hat{\Theta}_{\text{prop-2}}$ is not sensitive to μ_n at all. The performance of $\hat{\Theta}_{\text{RBLZ}}$ depends on the choice of μ_n . Therefore, it is difficult to make a fair comparison of $\hat{\Theta}_{\text{prop-2}}$ and $\hat{\Theta}_{\text{RBLZ}}$ under the same amount of L_1 regularization for scheme I. In Table 4, the condition number of the true precision matrix and the averaged condition numbers of $\hat{\Theta}_{\text{prop-2}}$ and $\hat{\Theta}_{\text{RBLZ}}^*$ are presented for schemes II and III only. Table 4 reveals that the averaged condition numbers of $\hat{\Theta}_{\text{RBLZ}}^*$ are larger than the true values and have relatively large standard errors.

6. Data Application

To illustrate the applicability of the proposed method, we applied the precision matrix estimator to call center data (available at <http://iew3.technion.ac.il/serveng2012S/callcenterdata/index.html>). The data recorded the time of the phone calls entering the call center of “Anonymous Bank” in Israel every day in 1999. Because of the difference of the arrival patterns between the weekdays and weekends, we discarded the data for the weekends (Friday and

Table 4. Comparison of the averaged condition numbers of $\widehat{\Theta}_{\text{prop-2}}$ and $\widehat{\Theta}_{\text{RBLZ}}^*$ for schemes II and III. Results are averaged over 400 replications with the standard errors shown in the brackets.

Scheme	p_n	$\text{cond}(\Theta_0)$	$\text{cond}(\widehat{\Theta}_{\text{prop-2}})$	$\text{cond}(\widehat{\Theta}_{\text{RBLZ}}^*)$
II	100	453.34	309.51 (4.59)	491.00 (29.70)
	200	453.34	344.18 (4.43)	645.07 (22.85)
	400	453.34	383.53 (4.42)	914.96 (47.99)
III	100	992.11	741.21 (9.50)	1,247.41 (41.03)
	200	1,127.59	862.59 (9.17)	1,710.53 (71.18)
	400	1,176.39	982.94 (9.23)	2,360.38 (132.83)

Saturday), and only used those on the weekdays (258 days). Since there are relatively fewer calls before 7:00am, we only considered the time period from 7:00am to midnight. On each day, we divided the 17-hour period into 3-minute intervals and counted the number of calls $X_{i,j}$ for the i th day and j th time period with $i \in \{1, \dots, 258\}$ and $j \in \{1, \dots, 340\}$.

We aimed to use the arrival counts in the first half of the day to predict those in the second half of the day. Take $\mathbf{X}_i = (X_{i,1}, \dots, X_{i,340})^T$ to be the vector of observations on the i th day, for $i = 1, \dots, 258$, and let $\mathbf{X}_i^{(1)} = (X_{i,1}, \dots, X_{i,170})^T$ and $\mathbf{X}_i^{(2)} = (X_{i,171}, \dots, X_{i,340})^T$ be the observations in the first and second halves of the day, respectively. We partitioned the mean vector and covariance matrix of \mathbf{X}_i correspondingly by

$$\boldsymbol{\mu}_0 = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \boldsymbol{\Sigma}_0 = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}.$$

The best linear predictor of $\mathbf{X}_i^{(2)}$ is expressed as

$$\widehat{\mathbf{X}}_i^{(2)} = \boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}(\mathbf{X}_i^{(1)} - \boldsymbol{\mu}_1). \tag{6.1}$$

We used $\{\mathbf{X}_1, \dots, \mathbf{X}_{100}\}$ as the training set and $\{\mathbf{X}_{101}, \dots, \mathbf{X}_{258}\}$ as the testing set. The estimates of $\boldsymbol{\mu}_0$ and $\boldsymbol{\Sigma}_0$, calculated based on the training data, were plugged into (6.1) for prediction. The sample mean $100^{-1} \sum_{i=1}^{100} \mathbf{X}_i$ was used to estimate $\boldsymbol{\mu}_0$ while the inverses of different precision matrix estimates were applied to estimate $\boldsymbol{\Sigma}_0$. Specifically, we calculated $\widehat{\Theta}_{\text{prop-1}}$, $\widehat{\Theta}_{\text{prop-2}}$, $\widehat{\Theta}_{\text{WLKR}}$, $\widehat{\Theta}_{\text{RBLZ}}$, $\widehat{\Theta}_{\text{GLasso}}$, and $\widehat{\Theta}_{\text{banded}}$ as the estimates of Θ_0 first, and took their inverses to estimate $\boldsymbol{\Sigma}_0$. For the precision matrix estimation problem, $(n, p_n) = (100, 340)$ and the tuning parameters were selected in the same way as that in Section 5 with grid points $\kappa_n \in \{1.2143^i : i = 0, 1, \dots, 19\}$ and $\mu_n \in \{0.1 \times 1.1708^i : i = 0, 1, \dots, 19\}$. The performance of different methods for estimating Θ_0 was compared by the averaged absolute forecast error (AFE) based on the testing

Table 5. Comparison of the AFE based on $\hat{\Theta}_{\text{prop-1}}$, $\hat{\Theta}_{\text{prop-2}}$, $\hat{\Theta}_{\text{WLKR}}$, $\hat{\Theta}_{\text{RBLZ}}$, $\hat{\Theta}_{\text{GLasso}}$, and $\hat{\Theta}_{\text{banded}}$ for the call center data, with SE denoting the standard error.

$\hat{\Theta}$	$\hat{\Theta}_{\text{prop-1}}$	$\hat{\Theta}_{\text{prop-2}}$	$\hat{\Theta}_{\text{WLKR}}$	$\hat{\Theta}_{\text{RBLZ}}$	$\hat{\Theta}_{\text{GLasso}}$	$\hat{\Theta}_{\text{banded}}$
AFE	1.83	1.76	1.83	1.81	1.84	1.82
SE	0.01	0.01	0.01	0.01	0.01	0.01

data. Specifically,

$$\text{AFE} = \frac{1}{158 \times 170} \sum_{i=101}^{258} \sum_{j=1}^{170} |e_{j,170}^T (\widehat{\mathbf{X}}_i^{(2)} - \mathbf{X}_i^{(2)})|.$$

Table 5 presents the AFE of the best linear predictor calculated with Θ_0 estimated by different methods. From Table 5, it is clear that $\hat{\Theta}_{\text{prop-2}}$ corresponds to a smaller averaged absolute forecast error than the other estimators.

Supplementary Materials

The detailed proofs of Theorems 1–4 and Proposition 1 are relegated to the Supplementary Material.

Acknowledgment

C. Zhang’s research is supported by the NSF grants DMS-1308872 and DMS-1521761, Wisconsin Alumni Research Foundation, and National Natural Science Foundation of China, grants 11690014. X. Guo’s research is supported by National Natural Science Foundation of China, grants 11601500, and the Fundamental Research Funds for the Central Universities.

References

- Bai, Z. D. and Yin, Y. Q. (1993). Limit of the smallest eigenvalue of a large dimensional sample covariance matrix. *Ann. Prob.* **21**, 1275–1294.
- Banerjee, O., El Ghaoui, L. and d’Aspremont, A. (2008). Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data. *J. Mach. Learn. Res.* **9**, 485–516.
- Bickel, P. J. and Levina, E. (2008). Regularized estimation of large covariance matrices. *Ann. Statist.* **36**, 199–227.
- Boyd, S., Parikh, N., Chu, E., Peleato, B. and Eckstein, J. (2010). Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundat. Trends Mach. Learn.* **3**, 1–122.
- Cai, T. T. and Liu, W. (2011). Adaptive thresholding for sparse covariance matrix estimation. *J. Amer. Statist. Assoc.* **106**, 672–684.

- Cai, T. T. and Zhou, H. H. (2012). Minimax estimation of large covariance matrices under ℓ_1 -norm. *Statistica Sinica* **22**, 1319–1378.
- DeGroot, M. H. (2004). *Optimal Statistical Decisions*. Wiley-Interscience.
- El Karoui, N. (2009). Concentration of measure and spectra of random matrices: applications to correlation matrices, elliptical distributions and beyond. *Ann. Appl. Prob.* **19**, 2362–2405.
- Farrell, R. H. (1985). *Multivariate Calculation*. Springer, New York.
- Friedman, J., Hastie, T. and Tibshirani, R. (2008). Sparse inverse covariance estimation with the graphical lasso. *Biostatistics* **9**, 432–441.
- Lam, C. and Fan, J. (2009). Sparsistency and rates of convergence in large covariance matrix estimation. *Ann. Statist.* **37**, 4254–4278.
- Liu, H., Wang, L. and Zhao, T. (2014). Sparse covariance matrix estimation with eigenvalue constraints. *J. Comput. Graph. Statist.* **23**, 439–459.
- Marčenko, V. A. and Pastur, L. A. (1967). Distribution of eigenvalues for some sets of random matrices. *Math. USSR-Sb.* **1**, 507–536.
- Meinshausen, N. and Bühlmann, P. (2006). High-dimensional graphs and variable selection with the lasso. *Ann. Statist.* **34**, 1436–1462.
- Peng, J., Wang, P., Zhou, N. and Zhu, J. (2009). Partial correlation estimation by joint sparse regression models. *J. Amer. Statist. Assoc.* **104**, 735–746.
- Ravikumar, P., Wainwright, M. J., Raskutti, G. and Yu, B. (2011). High-dimensional covariance estimation by minimizing ℓ_1 -penalized log-determinant divergence. *Electron. J. Statist.* **5**, 935–980.
- Rothman, A. J. (2012). Positive definite estimators of large covariance matrices. *Biometrika* **99**, 733–740.
- Rothman, A. J., Bickel, P. J., Levina, E. and Zhu, J. (2008). Sparse permutation invariant covariance estimation. *Electron. J. Statist.* **2**, 494–515.
- Rothman, A. J., Levina, E. and Zhu, J. (2009). Generalized thresholding of large covariance matrices. *J. Amer. Statist. Assoc.* **104**, 177–186.
- Silverstein, J. W. and Choi, S. I. (1995). Analysis of the limiting spectral distribution of large dimensional random matrices. *J. Mult. Anal.* **54**, 295–309.
- Stewart, G. W. and Sun, J. (1990). *Matrix Perturbation Theory*. Academic Press, Boston.
- Won, J., Lim, J., Kim, S. and Rajaratnam, B. (2013). Condition-number-regularized covariance estimation. *J. R. Statist. Soc. B* **75**, 427–450.
- Xue, L., Ma, S. and Zou, H. (2012). Positive-definite ℓ_1 -penalized estimation of large covariance matrices. *J. Amer. Statist. Assoc.* **107**, 1480–1491.

Department of Statistics and Finance, School of Management, University of Science and Technology of China, Hefei, Anhui 230026, China.

E-mail: xiaoguo@ustc.edu.cn

School of Mathematical Sciences, Nankai University, Tianjin 300071, China;

Department of Statistics, University of Wisconsin, Madison, WI 53706, USA.

E-mail: cmzhang@stat.wisc.edu

(Received January 2015; accepted June 2016)

The Effect of L_1 Penalization on Condition Number Constrained Estimation of Precision Matrix

Xiao Guo and Chunming Zhang

University of Science and Technology of China and University of Wisconsin-Madison

Supplementary Material

The supplementary material here includes the detailed proofs of Theorems 1–4 and Proposition 1 in the paper.

S1 Detailed proofs

Proof of Theorem 1. Without loss of generality, we assume $\boldsymbol{\mu}_0 = \mathbf{0}$.

Condition A1 indicates that $P\{|\mathbf{S}_n(i, j) - \boldsymbol{\Sigma}_0(i, j)| \geq \delta\} \leq C \exp(-c\delta^2 n)$ for $i, j = 1, \dots, p_n$ with an arbitrarily small constant $\delta \in (0, \infty)$ (see, for example, (11) and Lemma A.3 of Bickel and Levina (2008)), and hence

$$\|\widehat{W}^2 - W_0^2\|_2^2 = O_P\{\log(p_n)/n\} = \|\widehat{W}^{-1} - W_0^{-1}\|_2^2. \quad (\text{S1.1})$$

Under Condition A2, we have $P\{|\mathbf{S}_n(i, j) - \boldsymbol{\Sigma}_0(i, j)| \geq \delta\} \leq Cn^{-\beta/4}\delta^{-\beta/2}$ for $i, j = 1, \dots, p_n$ with a constant $\delta \in (0, \infty)$ (see, for example, Lemma 2 of Ravikumar et al. (2011)), which implies

$$\|\widehat{W}^2 - W_0^2\|_2^2 = O_P(p_n^{4/\beta}/n) = \|\widehat{W}^{-1} - W_0^{-1}\|_2^2. \quad (\text{S1.2})$$

It's easy to see that Condition A3 implies

$$\|\widehat{W}^2 - W_0^2\|_2^2 = O_P(p_n/n) = \|\widehat{W}^{-1} - W_0^{-1}\|_2^2.$$

Therefore, under either Condition A1 or A2 or A3,

$$\|\widehat{W}^2 - W_0^2\|_2^2 = o_P(1) = \|\widehat{W}^{-1} - W_0^{-1}\|_2^2.$$

To prove $\|\widehat{\boldsymbol{\Theta}}_{\text{prop-1}} - \boldsymbol{\Theta}_0\|_2 \xrightarrow{P} 0$, it suffices to show that $\|\widetilde{\Omega}_{\kappa_n} - \Omega_0\|_2 \xrightarrow{P} 0$.

Under Condition A1 or A2, $\boldsymbol{\Sigma}_0$ being diagonal induces $\Gamma_0 = \mathbf{I}_{p_n}$. From (2.2), $\widetilde{\Omega}_{\kappa_n} = \{p_n/\text{tr}(\mathbf{R}_n)\}\mathbf{I}_{p_n} = \mathbf{I}_{p_n}$ due to $\kappa_n = 1$. Hence, $\|\widetilde{\Omega}_{\kappa_n} - \Omega_0\|_2 \xrightarrow{P} 0$.

Under Condition A3, we first prove $\|\mathbf{S}_n - \boldsymbol{\Sigma}_0\|_2 \xrightarrow{P} 0$ for $\boldsymbol{\Sigma}_0 = \mathbf{I}_{p_n}$. For $i = 1, \dots, n$, define $\mathbf{X}_i^* = (\mathbf{X}_i^T, \mathbf{Y}_i^T)^T \in \mathbb{R}^{p_n^*}$ with $p_n^* > p_n$ an integer and $\mathbf{Y}_1, \dots, \mathbf{Y}_n \in \mathbb{R}^{p_n^* - p_n}$ i.i.d. random vectors, such that $\{\mathbf{e}_{j, p_n^*}^T \mathbf{X}_i^* : i = 1, \dots, n; j = 1, \dots, p_n^*\}$ are i.i.d. random variables. Let $\mathbf{S}_n^* = n^{-1} \sum_{i=1}^n \mathbf{X}_i^* \mathbf{X}_i^{*T}$.

From Theorem 2 of Bai and Yin (1993), if $\lim_{n \rightarrow \infty} p_n^*/n = y$ with a constant $y \in (0, 1)$, then $\lambda_{\max}(\mathbf{S}_n^*) \xrightarrow{P} (1 + \sqrt{y})^2$ and $\lambda_{\min}(\mathbf{S}_n^*) \xrightarrow{P} (1 - \sqrt{y})^2$. We know $\lambda_{\min}(\mathbf{S}_n^*) \leq \lambda_{\min}(n^{-1} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T) \leq \lambda_{\max}(n^{-1} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T) \leq \lambda_{\max}(\mathbf{S}_n^*)$. Thus, if y is arbitrarily close to 0, then we have $\lambda_{\max}(n^{-1} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T) \xrightarrow{P} 1$ and $\lambda_{\min}(n^{-1} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T) \xrightarrow{P} 1$. From

$$\|\mathbf{S}_n - \boldsymbol{\Sigma}_0\|_2 \leq \left\| n^{-1} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T - \boldsymbol{\Sigma}_0 \right\|_2 + \|\bar{\mathbf{X}} \bar{\mathbf{X}}^T\|_2 = \text{I} + \text{II},$$

$\text{I} = \max\{|\lambda_{\max}(n^{-1} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T) - 1|, |\lambda_{\min}(n^{-1} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^T) - 1|\} \xrightarrow{P} 0$ and $\text{II} = \|\bar{\mathbf{X}}^T \bar{\mathbf{X}}\|_2 \xrightarrow{P} 0$, we have $\|\mathbf{S}_n - \boldsymbol{\Sigma}_0\|_2 \xrightarrow{P} 0$.

For $\boldsymbol{\Sigma}_0$ not necessarily equal to \mathbf{I}_{p_n} , $\|\mathbf{S}_n - \boldsymbol{\Sigma}_0\|_2 \leq \|\boldsymbol{\Sigma}_0^{1/2}\|_2 \|\boldsymbol{\Sigma}_0^{-1/2} \mathbf{S}_n \boldsymbol{\Sigma}_0^{-1/2} - \mathbf{I}_{p_n}\|_2 \|\boldsymbol{\Sigma}_0^{1/2}\|_2 \xrightarrow{P} 0$, since $\boldsymbol{\Sigma}_0^{-1/2} \mathbf{S}_n \boldsymbol{\Sigma}_0^{-1/2}$ is the sample covariance matrix of $\{\boldsymbol{\Sigma}_0^{-1/2} \mathbf{X}_1, \dots, \boldsymbol{\Sigma}_0^{-1/2} \mathbf{X}_n\}$ which are i.i.d. with covariance matrix \mathbf{I}_{p_n} and

$$\begin{aligned} \{E(|\mathbf{e}_{1, p_n}^T \boldsymbol{\Sigma}_0^{-1/2} \mathbf{X}_1|^4)\}^{1/4} &\leq \sum_{i=1}^{p_n} \{E(|\mathbf{e}_{1, p_n}^T \boldsymbol{\Sigma}_0^{-1/2} \mathbf{e}_{i, p_n} X_{1, i}|^4)\}^{1/4} \\ &\leq \max_{1 \leq i \leq p_n} \{E(|X_{1, i}|^4)\}^{1/4} \sum_{i=1}^{p_n} |\mathbf{e}_{1, p_n}^T \boldsymbol{\Sigma}_0^{-1/2} \mathbf{e}_{i, p_n}| = \max_{1 \leq i \leq p_n} \{E(|X_{1, i}|^4)\}^{1/4} \|\boldsymbol{\Sigma}_0^{-1/2}\|_\infty \\ &< C < \infty. \end{aligned}$$

Therefore, $\|\mathbf{R}_n - \Gamma_0\|_2 \xrightarrow{P} 0$, which implies that $\|\tilde{\Omega}_{\kappa_n} - \Omega_0\|_2 \xrightarrow{P} 0$ since $\liminf_{n \rightarrow \infty} \{\kappa_n - \text{cond}(\Omega_0)\} > 0$.

The result $\|\hat{\boldsymbol{\Theta}}_{\text{prop-1}}^{-1} - \boldsymbol{\Sigma}_0\|_2 \xrightarrow{P} 0$ comes from $\|\hat{\boldsymbol{\Theta}}_{\text{prop-1}}^{-1} - \boldsymbol{\Sigma}_0\|_2 = O_P(\|\hat{\boldsymbol{\Theta}}_{\text{prop-1}} - \boldsymbol{\Theta}_0\|_2)$. ■

Proof of Theorem 2. Suppose the eigendecomposition of \mathbf{R}_n is $Q \text{diag}(\hat{\lambda}_1, \dots, \hat{\lambda}_{p_n}) Q^T$, where $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_{p_n}$ are the eigenvalues of \mathbf{R}_n . From Won et al. (2013), $\tilde{\Omega}_{\kappa_n}^{-1} = Q \text{diag}(\tilde{\lambda}_1, \dots, \tilde{\lambda}_{p_n}) Q^T$, where $\tilde{\lambda}_i = \min\{\max(\tau^*, \hat{\lambda}_i), \kappa_n \tau^*\}$ with $\tau^* \in (0, \infty)$ depending on $\hat{\lambda}_1, \dots, \hat{\lambda}_{p_n}$ and κ_n . Hence, $\tilde{\Omega}_{\kappa_n}^{-1}$ truncates the eigenvalues of \mathbf{R}_n . From Stewart and Sun (1990) (Corollary 4.10, p. 203),

$$\max_{1 \leq i \leq p_n} |\tilde{\lambda}_i - \lambda_i| \leq \|\tilde{\Omega}_{\kappa_n}^{-1} - \Gamma_0\|_2,$$

where $\lambda_1 \geq \dots \geq \lambda_{p_n}$ are the eigenvalues of Γ_0 . If $\|\tilde{\Omega}_{\kappa_n}^{-1} - \Gamma_0\|_2 \xrightarrow{P} 0$, then $\max_{1 \leq i \leq p_n} |\tilde{\lambda}_i - \lambda_i| \xrightarrow{P} 0$ which implies that $\mathbb{F}^{\tilde{\Omega}_{\kappa_n}^{-1}}$ converges weakly to F_0 in probability. Therefore, to prove $\|\hat{\Theta}_{\text{prop-1}} - \Theta_0\|_2 \rightarrow 0$ in probability, we only need to show $\|\tilde{\Omega}_{\kappa_n}^{-1} - \Gamma_0\|_2 \rightarrow 0$ in probability, and it suffices to show that $\mathbb{F}^{\tilde{\Omega}_{\kappa_n}^{-1}}$ doesn't converge weakly to F_0 in probability.

Under Condition B1, if $\lim_{n \rightarrow \infty} p_n/n = \infty$, then the rank of \mathbf{R}_n is at most n when $p_n > n$, and hence, the proportion of the 0 eigenvalues among $\hat{\lambda}_1, \dots, \hat{\lambda}_{p_n}$ is at least $(p_n - n)/p_n$ which converges to 1 as $n \rightarrow \infty$. Therefore, $\mathbb{F}^{\mathbf{R}_n}$ will converge weakly to $I_{[0, \infty)}$ in probability. Since $\tilde{\lambda}_i = \min\{\max(\tau^*, \hat{\lambda}_i), \kappa_n \tau^*\}$, if $\mathbb{F}^{\tilde{\Omega}_{\kappa_n}^{-1}}$ converges weakly in probability, then the limit is $I_{[c, \infty)}$ for some $c \in [0, \infty)$. Since $F_0 \neq I_{[c, \infty)}$ for any $C \in [0, \infty)$, $\mathbb{F}^{\tilde{\Omega}_{\kappa_n}^{-1}}$ doesn't converge weakly to F_0 in probability. Therefore, $\|\hat{\Theta}_{\text{prop-1}} - \Theta_0\|_2 \rightarrow 0$ in probability.

Under Condition B2, we will show that $|\text{cond}(\tilde{\Omega}_{\kappa_n}^{-1}) - \text{cond}(\Gamma_0)| \rightarrow 0$ in probability which implies that $\|\tilde{\Omega}_{\kappa_n}^{-1} - \Gamma_0\|_2 \rightarrow 0$ in probability. From Theorem 1 in Won et al. (2013), $\text{cond}(\tilde{\Omega}_{\kappa_n}^{-1}) = \min\{\kappa_n, \text{cond}(\mathbf{R}_n)\}$. Since $|\min\{\kappa_n, \text{cond}(\mathbf{R}_n)\} - \text{cond}(\Gamma_0)| \rightarrow 0$ in probability, we have $|\text{cond}(\tilde{\Omega}_{\kappa_n}^{-1}) - \text{cond}(\Gamma_0)| \rightarrow 0$ in probability as $n \rightarrow \infty$.

Under Condition B3, we will show that $\mathbb{F}^{\tilde{\Omega}_{\kappa_n}^{-1}}$ does not converge weakly to F_0 in probability. We truncate $\mathbb{F}^{\mathbf{R}_n}$ in order to obtain $\mathbb{F}^{\tilde{\Omega}_{\kappa_n}^{-1}}$, i.e., $\mathbb{F}^{\tilde{\Omega}_{\kappa_n}^{-1}} = \mathbb{F}^{\mathbf{R}_n} I_{[\tau^*, \kappa_n \tau^*)} + I_{[\kappa_n \tau^*, \infty)}$. If $\mathbb{F}^{\tilde{\Omega}_{\kappa_n}^{-1}}$ converges weakly to F_0 in probability, then $F_0 = F I_{[l_{\min}, l_{\max})} + I_{[l_{\max}, \infty)}$, which contradicts Condition B3.

Therefore, we have demonstrated that $\|\hat{\Theta}_{\text{prop-1}} - \Theta_0\|_2 \rightarrow 0$ in probability under either Condition B1 or B2 or B3. Next, we will show that $\|\hat{\Theta}_{\text{prop-1}}^{-1} - \Sigma_0\|_2 \rightarrow 0$ in probability. If $\|\hat{\Theta}_{\text{prop-1}}^{-1} - \Sigma_0\|_2 = o_P(1)$, then $\|\hat{\Theta}_{\text{prop-1}} - \Theta_0\|_2 = o_P(1)$ because $\|\hat{\Theta}_{\text{prop-1}} - \Theta_0\|_2 = O_P(\|\hat{\Theta}_{\text{prop-1}}^{-1} - \Sigma_0\|_2)$. Since $\|\hat{\Theta}_{\text{prop-1}} - \Theta_0\|_2 \rightarrow 0$ in probability, we can claim $\|\hat{\Theta}_{\text{prop-1}}^{-1} - \Sigma_0\|_2 \rightarrow 0$ in probability. ■

Proof of Theorem 3. Following the proofs of Corollaries 1 and 2 in Ravikumar et al. (2011), we have that, with probability tending to 1,

$$|\hat{\Omega}_{\text{RBLZ}} - \Omega_0|_{\infty}^2 \leq C r_n^*$$

with $r_n^* = \log(p_n)/n$ under Condition C1, and $r_n^* = p_n^{4\tau/\beta}/n$ under Condition C2, where $|\cdot|_{\infty}$ is the matrix elementwise L_{∞} norm defined as $|A|_{\infty} = \max_{i,j} |A(i, j)|$ for a generic matrix A . The proof of Theorem 1 in Ravikumar et al. (2011) indicates that $\{(i, j) :$

$\widehat{\Omega}_{\text{RBLZ}}(i, j) \neq 0\} \subseteq \{(i, j) : \Omega_0(i, j) \neq 0\}$ with probability tending to 1.

For $n = 1, 2, \dots$, let \mathcal{A}_n denote the event that $|\widehat{\Omega}_{\text{RBLZ}} - \Omega_0|_\infty^2 \leq Cr_n^*$ and $\{(i, j) : \widehat{\Omega}_{\text{RBLZ}}(i, j) \neq 0\} \subseteq \{(i, j) : \Omega_0(i, j) \neq 0\}$. Hence, $\lim_{n \rightarrow \infty} P(\mathcal{A}_n) = 1$. Then, conditional on event \mathcal{A}_n , we have

$$\begin{aligned} & \|\widehat{\Omega}_{\text{RBLZ}} - \Omega_0\|_2^2 \leq \|\widehat{\Omega}_{\text{RBLZ}} - \Omega_0\|_F^2 \\ &= \sum_{i=1}^{p_n} |\widehat{\Omega}_{\text{RBLZ}}(i, i) - \Omega_0(i, i)|^2 + \sum_{i \neq j: \Omega_0(i, j) \neq 0} |\widehat{\Omega}_{\text{RBLZ}}(i, j) - \Omega_0(i, j)|^2 \quad (\text{S1.3}) \end{aligned}$$

and

$$\|\widehat{\Omega}_{\text{RBLZ}} - \Omega_0\|_2^2 \leq \|\widehat{\Omega}_{\text{RBLZ}} - \Omega_0\|_\infty^2 \leq t_n^2 |\widehat{\Omega}_{\text{RBLZ}} - \Omega_0|_\infty^2. \quad (\text{S1.4})$$

If $p_n \leq s_n$, then (S1.3) and (S1.4) indicate $\|\widehat{\Omega}_{\text{RBLZ}} - \Omega_0\|_2^2 \leq C \min(p_n + s_n, t_n^2) r_n^* \leq Cr_n$ under Condition C1 or C2. Next, we consider the case $p_n > s_n$. Define $t_n^i = |\{j = 1, \dots, p_n : \Theta_0(i, j) \neq 0\}|$. For any $i \in \{1, \dots, p_n\}$ such that $t_n^i = 1$, we know $\Omega_0 \mathbf{e}_{i, p_n} = \mathbf{e}_{i, p_n}$, which means that the diagonal element is the only nonzero element in the i th column of Ω_0 . Since $p_n > s_n$, we have $|\{i = 1, \dots, p_n : t_n^i = 1\}| \geq p_n - s_n$. Because $\{(i, j) : \widehat{\Omega}_{\text{RBLZ}}(i, j) \neq 0\} \subseteq \{(i, j) : \Omega_0(i, j) \neq 0\}$, from the definition of $\widehat{\Omega}_{\text{RBLZ}}$, we have $\widehat{\Omega}_{\text{RBLZ}} \mathbf{e}_{i, p_n} = \mathbf{e}_{i, p_n}$ for any $i \in \{1, \dots, p_n\}$ with $\Omega_0 \mathbf{e}_{i, p_n} = \mathbf{e}_{i, p_n}$. Hence, $|\widehat{\Omega}_{\text{RBLZ}}(i, i) - \Omega_0(i, i)| = 0$ for $i \in \{1, \dots, p_n\}$ with $t_n^i = 1$. Therefore, (S1.3) indicates $\|\widehat{\Omega}_{\text{RBLZ}} - \Omega_0\|_2^2 \leq C(1 + s_n) r_n^*$, which together with (S1.4) implies that $\|\widehat{\Omega}_{\text{RBLZ}} - \Omega_0\|_2^2 \leq C \min(1 + s_n, t_n^2) r_n^* \leq Cr_n$.

Hence, under Condition C1 or C2, $\|\widehat{\Omega}_{\text{RBLZ}} - \Omega_0\|_2^2 = O_P(r_n)$. Therefore, from (S1.1) and (S1.2),

$$\begin{aligned} & \|\widehat{\Theta}_{\text{RBLZ}} - \Theta_0\|_2 = \|\widehat{W}^{-1} \widehat{\Omega}_{\text{RBLZ}} \widehat{W}^{-1} - W_0^{-1} \Omega_0 W_0^{-1}\|_2 \\ & \leq \|\widehat{W}^{-1} - W_0^{-1}\|_2 \|\widehat{\Omega}_{\text{RBLZ}} - \Omega_0\|_2 \|\widehat{W}^{-1} - W_0^{-1}\|_2 \\ & \quad + \|\widehat{W}^{-1} - W_0^{-1}\|_2 (\|\widehat{\Omega}_{\text{RBLZ}}\|_2 \|W_0^{-1}\|_2 + \|\widehat{W}^{-1}\|_2 \|\Omega_0\|_2) \\ & \quad + \|\widehat{\Omega}_{\text{RBLZ}} - \Omega_0\|_2 \|\widehat{W}^{-1}\|_2 \|W_0^{-1}\|_2 = O_P(r_n^{1/2}). \end{aligned}$$

We obtain $\|\widehat{\Theta}_{\text{RBLZ}}^{-1} - \Sigma_0\|_2^2 = O_P(r_n)$, since $\|\widehat{\Theta}_{\text{RBLZ}}^{-1} - \Sigma_0\|_2^2 = O_P(\|\widehat{\Theta}_{\text{RBLZ}} - \Theta_0\|_2^2)$. ■

Proof of Theorem 4. Following the proof of Theorem 3, under Condition C1 or C2, $\|\widehat{\Omega}_{\text{RBLZ}} - \Omega_0\|_2^2 = O_P(r_n)$. Now that $\text{cond}(\widehat{\Omega}_{\text{RBLZ}}) - \text{cond}(\Omega_0) = o_P(1)$, from

$\liminf_{n \rightarrow \infty} \{\kappa_n - \text{cond}(\Omega_0)\} > 0$, we have $\text{cond}(\widehat{\Omega}_{\text{RBLZ}}) \leq \kappa_n$ with probability tending to 1, which means that $\widehat{\Omega}_{\mu_n, \kappa_n} = \widehat{\Omega}_{\text{RBLZ}}$ with probability tending to 1, and hence $\lim_{n \rightarrow \infty} \text{P}(\widehat{\Theta}_{\text{prop-2}} = \widehat{\Theta}_{\text{RBLZ}}) = 1$. Therefore, from the conclusion in Theorem 3, $\|\widehat{\Theta}_{\text{prop-2}} - \Theta_0\|_2^2 = O_{\text{P}}(r_n) = \|\widehat{\Theta}_{\text{prop-2}}^{-1} - \Sigma_0\|_2^2$. ■

Proof of Proposition 1. From (3.2) and (3.3), suppose the eigendecomposition of variable Ω is $RM R^T$, where R is orthogonal and $M = \text{diag}(m_1, \dots, m_{p_n})$ with $m_1 \leq \dots \leq m_{p_n}$. For Step 1 in Section 3,

$$\begin{aligned}
 & \arg \min_{\Omega > 0, \text{cond}(\Omega) \leq \kappa_n} L_\rho(\Omega, Z^{(i-1)}; U^{(i-1)}) \\
 = & \arg \min_{\Omega > 0, \text{cond}(\Omega) \leq \kappa_n} -\log\{\det(\Omega)\} + \text{tr}(\mathbf{R}_n \Omega) + \frac{\rho}{2} \|\Omega - Z^{(i-1)} + U^{(i-1)}\|_F^2 \\
 = & \arg \min_{\Omega > 0, \text{cond}(\Omega) \leq \kappa_n} -\log\{\det(\Omega)\} + \text{tr}(\mathbf{R}_n \Omega) + \frac{\rho}{2} \text{tr}\{\Omega \Omega^T + 2(-Z^{(i-1)} + U^{(i-1)})\Omega^T\} \\
 = & \arg \min_{\Omega > 0, \text{cond}(\Omega) \leq \kappa_n} -\log\{\det(\Omega)\} + \frac{\rho}{2} \text{tr}(\Omega \Omega^T) + \rho \text{tr}\{(\mathbf{R}_n / \rho - Z^{(i-1)} + U^{(i-1)})\Omega^T\} \\
 = & \arg \min_{\Omega > 0, \text{cond}(\Omega) \leq \kappa_n} -\log\{\det(\Omega)\} + \frac{\rho}{2} \text{tr}(\Omega \Omega^T) + \rho \text{tr}\{(VDV^T)\Omega^T\} \\
 = & \arg \min_{\Omega = RM R^T: M > 0, \text{cond}(M) \leq \kappa_n} -\log\{\det(M)\} + \frac{\rho}{2} \text{tr}(MM^T) + \rho \text{tr}\{(VDV^T)(RM R^T)^T\} \\
 = & \arg \min_{\Omega = RM R^T: R=V, M > 0, \text{cond}(M) \leq \kappa_n} -\log\{\det(M)\} + \frac{\rho}{2} \text{tr}(MM^T) + \rho \text{tr}(DM^T). \quad (\text{S1.5})
 \end{aligned}$$

The last equation in (S1.5) is true since $\text{tr}\{(VDV^T)(RM R^T)^T\} \geq \text{tr}(DM^T)$ with equality if $R = V$ (Theorem 14.3.2 in Farrell (1985)). Therefore, to prove $\Omega^{(i)} = V\tilde{D}V^T$, it suffices to show that

$$\tilde{D} = \arg \min_{M: M > 0, \text{cond}(M) \leq \kappa_n} -\log\{\det(M)\} + \frac{\rho}{2} \text{tr}(MM^T) + \rho \text{tr}(DM^T),$$

which is equivalent to

$$\begin{aligned}
 \tilde{D} &= \arg \min_{M: 0 < m_1 \leq \dots \leq m_{p_n}, m_{p_n}/m_1 \leq \kappa_n} \left\{ -\sum_{j=1}^{p_n} \log(m_j) + \frac{\rho}{2} \sum_{j=1}^{p_n} m_j^2 + \rho \sum_{j=1}^{p_n} d_j m_j \right\} \\
 &= \arg \min_{M: \exists \tau, 0 < \tau \leq m_1 \leq \dots \leq m_{p_n} \leq \kappa_n \tau} \sum_{j=1}^{p_n} \left\{ -\log(m_j) + \frac{\rho}{2} (m_j + d_j)^2 \right\}. \quad (\text{S1.6})
 \end{aligned}$$

Define

$$g(m_j; d_j) = -\log(m_j) + \frac{\rho}{2} (m_j + d_j)^2.$$

Then, $g(m_j; d_j)$ is strictly convex in $m_j \in (0, \infty)$ for any $j = 1, \dots, p_n$, and has a unique minimizer $\delta_j = -d_j/2 + \sqrt{d_j^2/4 + 1/\rho}$. Noting that $0 < \delta_1 \leq \dots \leq \delta_{p_n}$, if

$\delta_{p_n}/\delta_1 \leq \kappa_n$, then $\tilde{D} = \text{diag}(\delta_1, \dots, \delta_{p_n})$ coincides with the solution to problem (S1.6) with any $\tau \in [\delta_{p_n}/\kappa_n, \delta_1]$.

For case $\delta_{p_n}/\delta_1 > \kappa_n$, we first consider minimizing the objective function in (S1.6) with respect to m_1, \dots, m_{p_n} separately. For any $\tau > 0$ and $j = 1, \dots, p_n$, it follows that

$$\begin{aligned} m_j^*(\tau) &:= \arg \min_{\tau \leq m_j \leq \kappa_n \tau} \sum_{k=1}^{p_n} g(m_k; d_k) = \arg \min_{\tau \leq m_j \leq \kappa_n \tau} g(m_j; d_j) = \min\{\max(\tau, \delta_j), \kappa_n \tau\} \\ &= \begin{cases} \tau, & \text{if } \delta_j < \tau, \\ \delta_j, & \text{if } \tau \leq \delta_j \leq \kappa_n \tau, \\ \kappa_n \tau, & \text{if } \delta_j > \kappa_n \tau. \end{cases} \end{aligned}$$

Since $\tau \leq m_1^*(\tau) \leq \dots \leq m_{p_n}^*(\tau) \leq \kappa_n \tau$ for any $\tau > 0$, problem (S1.6) amounts to

$$\arg \min_{M: \exists \tau > 0, m_j = m_j^*(\tau)} \sum_{j=1}^{p_n} g(m_j; d_j) = \arg \min_{M: \exists \tau > 0, m_j = m_j^*(\tau)} \sum_{j=1}^{p_n} g(m_j^*(\tau); d_j).$$

Therefore, to prove that \tilde{D} is the solution to the optimization problem in (S1.6), we only need to show that τ_0 is the minimizer of

$$f(\tau) := \sum_{j=1}^{p_n} g(m_j^*(\tau); d_j) = \sum_{j: \delta_j < \tau} g(\tau; d_j) + \sum_{j: \tau \leq \delta_j \leq \kappa_n \tau} g(\delta_j; d_j) + \sum_{j: \delta_j > \kappa_n \tau} g(\kappa_n \tau; d_j).$$

We can verify that $g(m_j^*(\tau); d_j)$ is a convex function of $\tau \in (0, \infty)$ and has a continuous first-order derivative with respect to $\tau \in (0, \infty)$, for any $j = 1, \dots, p_n$. Therefore, $f(\tau)$ is convex and continuously differentiable for $\tau \in (0, \infty)$. For $\alpha \in \{1, \dots, p_n\}$ and $\beta \in \{1, \dots, p_n\}$ such that $\beta - 1 \geq \alpha$, define

$$\begin{aligned} R_{\alpha, \beta} &= \{\tau : \delta_\alpha < \tau \leq \delta_{\alpha+1} \text{ and } \delta_{\beta-1} \leq \kappa_n \tau < \delta_\beta\}, \\ f_{\alpha, \beta}(\tau) &= \sum_{j=1}^{\alpha} g(\tau; d_j) + \sum_{j=\alpha+1}^{\beta-1} g(\delta_j; d_j) + \sum_{j=\beta}^{p_n} g(\kappa_n \tau; d_j). \end{aligned}$$

Then, $f(\tau) = f_{\alpha, \beta}(\tau)$ for $\tau \in R_{\alpha, \beta}$. Since $f''_{\alpha, \beta}(\tau) > 0$ for $\tau \in R_{\alpha, \beta}$, we know $f'(\tau)$ is strictly monotone increasing on $[\delta_1, \delta_{p_n}/\kappa_n]$. It's also easy to see that $f(\tau)$ is decreasing for $\tau \in (0, \delta_1]$ and increasing for $\tau \in [\delta_{p_n}/\kappa_n, \infty)$. Then, the unique minimizer of $f(\tau)$ is the value of $\tau \in [\delta_1, \delta_{p_n}/\kappa_n]$ such that $f'(\tau) = 0$.

The solution to $f'_{\alpha, \beta}(\tau) = 0$ for $\tau \in (0, \infty)$ is

$$\tau_{\alpha, \beta} = \left[-\rho \left(\sum_{j=1}^{\alpha} d_j + \kappa_n \sum_{j=\beta}^{p_n} d_j \right) + \left\{ \rho^2 \left(\sum_{j=1}^{\alpha} d_j + \kappa_n \sum_{j=\beta}^{p_n} d_j \right)^2 + 4\rho(\alpha + \kappa_n^2 p_n) \right\}^{1/2} \right] / 2$$

$$-\kappa_n^2\beta + \kappa_n^2)(\alpha + p_n - \beta + 1)\}^{1/2}] / \{2\rho(\alpha + \kappa_n^2 p_n - \kappa_n^2\beta + \kappa_n^2)\}.$$

Then, $\tau_{\alpha,\beta}$ is also the solution to $f'(\tau) = 0$ if and only if $\tau_{\alpha,\beta} \in R_{\alpha,\beta}$. This value of $\tau_{\alpha,\beta}$ is the same as τ_0 .

In practice, we can search over $\{R_{\alpha,\beta} : \alpha, \beta = 1, \dots, p_n\}$ to find α_0 and β_0 such that $\tau_{\alpha_0,\beta_0} \in R_{\alpha_0,\beta_0}$. Start the selection procedure from (α^*, β^*) , where $\alpha^* = 1$ and β^* is the smallest index in $\{1, \dots, p_n\}$ such that $\delta_{\beta^*} > \kappa_n \delta_{\alpha^*}$. If $\tau_{\alpha^*,\beta^*} \notin R_{\alpha^*,\beta^*}$, then move on to R_{α^*+1,β^*} , R_{α^*+1,β^*+1} or R_{α^*,β^*+1} for the selection of α_0 and β_0 . Specifically, if $\kappa_n \delta_{\alpha^*+1} < \delta_{\beta^*}$, then move on to R_{α^*+1,β^*} ; if $\kappa_n \delta_{\alpha^*+1} > \delta_{\beta^*}$, then go to R_{α^*,β^*+1} ; otherwise, continue searching α_0 and β_0 within R_{α^*+1,β^*+1} . Repeat the above procedure until condition $\tau_{\alpha,\beta} \in R_{\alpha,\beta}$ is satisfied. The procedure requires $O(p_n)$ operations. ■