

# Optimal Rates of Convergence for Covariance Matrix Estimation

T. Tony Cai<sup>1</sup>, Cun-Hui Zhang<sup>2</sup> and Harrison H. Zhou<sup>3</sup>

University of Pennsylvania, Rutgers University and Yale University

## Abstract

Covariance matrix plays a central role in multivariate statistical analysis. Significant advances have been made recently on developing both theory and methodology for estimating large covariance matrices. However, a minimax theory has yet been developed. In this paper we establish the optimal rates of convergence for estimating the covariance matrix under both the operator norm and Frobenius norm. It is shown that optimal procedures under the two norms are different and consequently matrix estimation under the operator norm is fundamentally different from vector estimation. The minimax upper bound is obtained by constructing a special class of tapering estimators and by studying their risk properties. A key step in obtaining the optimal rate of convergence is the derivation of the minimax lower bound. The technical analysis requires new ideas that are quite different from those used in the more conventional function/sequence estimation problems.

**Keywords:** Covariance matrix, Frobenius norm, minimax lower bound, operator norm, optimal rate of convergence, tapering.

**AMS 2000 Subject Classification:** Primary 62H12; secondary 62F12, 62G09.

---

<sup>1</sup>The research of Tony Cai was supported in part by NSF Grant DMS-0604954.

<sup>2</sup>The research of Cun-hui Zhang was supported in part by NSF Grants DMS 0504387 and DMS 0604571 and NSA Grant MDS 904-02-1-0063

<sup>3</sup>The research of Harrison Zhou was supported in part by NSF Career Award DMS-0645676.

# 1 Introduction

Suppose we observe independent and identically distributed  $p$ -variate random variables  $\mathbf{X}_1, \dots, \mathbf{X}_n$  with covariance matrix  $\Sigma_{p \times p}$  and the goal is to estimate the unknown matrix  $\Sigma_{p \times p}$  based on the sample  $\{\mathbf{X}_i : i = 1, \dots, n\}$ . This covariance matrix estimation problem is of fundamental importance in multivariate analysis. A wide range of statistical methodologies, including clustering analysis, principal component analysis, linear and quadratic discriminant analysis, regression analysis, require the estimation of the covariance matrices. With dramatic advances in technology, large high-dimensional data are now routinely collected in scientific investigations. Examples include climate studies, gene expression arrays, functional magnetic resonance imaging, risk management and portfolio allocation and web search problems. In such settings, the standard and most natural estimator, the sample covariance matrix, often performs poorly. See, for example, Muirhead (1987), Johnstone (2001), Bickel and Levina (2008a, b) and Fan, Fan and Lv (2007).

Regularization methods, originally developed in nonparametric function estimation, have recently been applied to estimate large covariance matrices. These include banding method in Wu and Pourahmadi (2009) and Bickel and Levina (2008a), tapering in Furrer and Bengtsson (2007), thresholding in Bickel and Levina (2008b) and El Karoui (2008), penalized estimation in Huang, Liu, Pourahmadi and Liu (2006), Lam and Fan (2007) and Rothman, Bickel, Levina and Zhu (2008), regularizing principal components in Johnstone and Lu (2004) and Zou, Hastie, and Tibshirani (2006). Asymptotic properties and convergence results have been given in several papers. In particular, Bickel and Levina (2008a, 2008b), El Karoui (2008) and Lam and Fan (2007) showed consistency of their estimators in operator norm and even obtained explicit rates of convergence. However, it is not clear whether any of these rates of convergence are optimal.

Despite recent progress on covariance matrix estimation there has been remarkably little fundamental theoretical study on optimal estimation. In this paper, we establish the optimal rate of convergence for estimating the covariance matrix as well as its inverse over a wide range of classes of covariance matrices. Both the operator norm and Frobenius norm are considered. It is shown that optimal procedures for these two norms are different and consequently matrix estimation under the operator norm is fundamentally different from vector estimation. In addition, the results also imply that the banding estimator given in Bickel and Levina (2008a) is sub-optimal under the operator norm and the performance can be significantly improved.

We begin by considering optimal estimation of the covariance matrix  $\Sigma$  over a class

of matrices that has been considered in Bickel and Levina (2008a). Both minimax lower and upper bounds are derived. We write  $a_n \asymp b_n$  if there are positive constants  $c$  and  $C$  independent of  $n$  such that  $c \leq a_n/b_n \leq C$ . For a matrix  $A$  its operator norm is defined as  $\|A\| = \sup_{\|x\|_2=1} \|Ax\|_2$ . We assume that  $p \leq \exp(\gamma n)$  for some constant  $\gamma > 0$ . Combining the results given in Section 3, we have the following optimal rate of convergence for estimating the covariance matrix under the operator norm.

**Theorem 1** *The minimax risk of estimating the covariance matrix  $\Sigma$  over the class  $\mathcal{P}_\alpha$  given in (3) satisfies*

$$\inf_{\hat{\Sigma}} \sup_{\mathcal{P}_\alpha} \mathbb{E} \left\| \hat{\Sigma} - \Sigma \right\|^2 \asymp \min \left\{ n^{-\frac{2\alpha}{2\alpha+1}} + \frac{\log p}{n}, \frac{p}{n} \right\}. \quad (1)$$

The minimax upper bound is obtained by constructing a class of tapering estimators and by studying their risk properties. It is shown that the estimator with the optimal choice of the tapering parameter attains the optimal rate of convergence. In comparison to some existing methods in the literature, the proposed procedure does not attempt to estimate each row/column optimally as a vector. In fact, our procedure does not optimally trade bias and variance for each row/column. As a vector estimator, it has larger variance than squared bias for each row/column. In other words, it is undersmoothed as a vector.

A key step in obtaining the optimal rate of convergence is the derivation of the minimax lower bound. The lower bound is established by using a testing argument, where at the core is a novel construction of a collection of least favorable multivariate normal distributions and the application of Assouad's Lemma and Le Cam's method. The technical analysis requires ideas that are quite different from those used in the more conventional function/sequence estimation problems.

In addition to the asymptotic analysis, we also carry out a small simulation study to investigate the finite sample performance of the proposed estimator. The tapering estimator is easy to implement. The numerical performance of the estimator is compared with that of the banding estimator introduced in Bickel and Levina (2008a). The simulation study shows that the proposed estimator has good numerical performance; it nearly uniformly outperforms the banding estimator.

The paper is organized as follows. In Section 2, after basic notations and definitions are introduced, we propose a tapering procedure for the covariance matrix estimation. Section 3 derives the optimal rate of convergence for estimation under the operator norm. The upper bound is obtained by studying the properties of the tapering estimators and the minimax lower bound is obtained by a testing argument. Section 4 considers optimal

estimation under the Frobenius norm. The problem of estimating the inverse of a covariance matrix is treated in Section 5. Section 6 investigates the numerical performance of our procedure by a simulation study. The technical proofs of auxiliary lemmas are given in Section 7.

## 2 Methodology

In this section we will introduce a tapering procedure for estimating the covariance matrix  $\Sigma_{p \times p}$  based on a random sample of  $p$ -variate observations  $\mathbf{X}_1, \dots, \mathbf{X}_n$ . The properties of the tapering estimators under the operator norm and Frobenius norm are then studied and used to establish the minimax upper bounds in Sections 3 and 4.

Given a random sample  $\{\mathbf{X}_1, \dots, \mathbf{X}_n\}$  from a population with covariance matrix  $\Sigma = \Sigma_{p \times p}$ , the sample covariance matrix is

$$\frac{1}{n-1} \sum_{l=1}^n (\mathbf{X}_l - \bar{\mathbf{X}}) (\mathbf{X}_l - \bar{\mathbf{X}})^T,$$

which is an unbiased estimate of  $\Sigma$ , and the maximum likelihood estimator of  $\Sigma$  is

$$\Sigma^* = (\sigma_{ij}^*)_{1 \leq i, j \leq p} = \frac{1}{n} \sum_{l=1}^n (\mathbf{X}_l - \bar{\mathbf{X}}) (\mathbf{X}_l - \bar{\mathbf{X}})^T \quad (2)$$

when  $\mathbf{X}_l$ 's are normally distributed. These two estimators are close to each other for large  $n$ . We shall construct estimators of the covariance matrix  $\Sigma$  by tapering the maximum likelihood estimator  $\Sigma^*$ .

Following Bickel and Levina (2008a) we consider estimating the covariance matrix  $\Sigma_{p \times p} = (\sigma_{ij})_{1 \leq i, j \leq p}$  over the following parameter space

$$\mathcal{F}_\alpha = \mathcal{F}_\alpha(M_0, M) = \left\{ \Sigma : \max_j \sum_i \{|\sigma_{ij}| : |i - j| > k\} \leq Mk^{-\alpha} \text{ for all } k, \text{ and } \lambda_{\max}(\Sigma) \leq M_0 \right\} \quad (3)$$

where  $\lambda_{\max}(\Sigma)$  is the maximum eigenvalue of the matrix  $\Sigma$ , and  $\alpha > 0$ ,  $M > 0$  and  $M_0 > 0$ . Note that the smallest eigenvalue of any covariance matrix in the parameter space  $\mathcal{F}_\alpha$  is allowed to be 0 which is more general than the assumption in equation (5) of Bickel and Levina (2008a). The parameter  $\alpha$  in (3), which essentially specifies the rate of decay for the covariances  $\sigma_{ij}$  as they move away from the diagonal, can be viewed as an analog of the smoothness parameter in nonparametric function estimation problems. The optimal rate of convergence for estimating  $\Sigma$  over the parameter space  $\mathcal{F}_\alpha(M_0, M)$  critically depends on the value of  $\alpha$ . Our estimators of the covariance matrix  $\Sigma$  are constructed by tapering the maximum likelihood estimator (2) as follows.

**Estimation Procedure** For a given even integer  $k$  with  $1 \leq k \leq p$ , we define a tapering estimator as

$$\hat{\Sigma} = \hat{\Sigma}_k = (w_{ij}\sigma_{ij}^*)_{p \times p} \quad (4)$$

where  $\sigma_{ij}^*$  are the entries in the maximum likelihood estimator  $\Sigma^*$  and the weights

$$w_{ij} = k_h^{-1} \{(k - |i - j|)_+ - (k_h - |i - j|)_+\} \quad (5)$$

where  $k_h = k/2$ . Without loss of generality we assume that  $k$  is even. Note that the weights  $w_{ij}$  can be rewritten as

$$w_{ij} = \begin{cases} 1 & \text{when } |i - j| \leq k_h \\ 2 - \frac{|i-j|}{k_h} & \text{when } k_h < |i - j| < k \\ 0 & \text{otherwise.} \end{cases} .$$

See Figure 1 for a plot of the weights  $w_{ij}$  as a function of  $|i - j|$ .

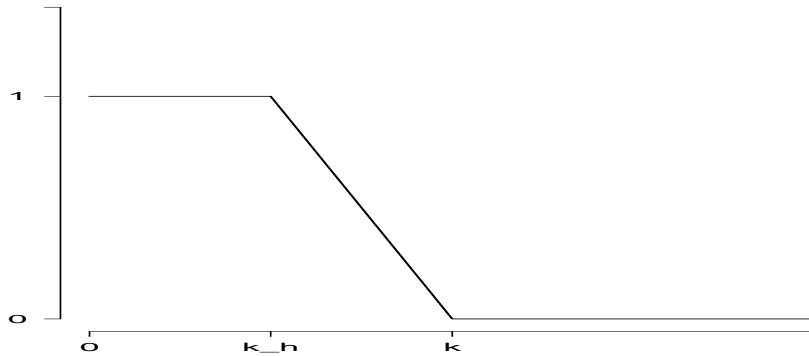


Figure 1: The weights as a function of  $|i - j|$ .

The tapering estimators are different from the banding estimators used in Bickel and Levina (2008a). It is important to note that the tapering estimator given in (4) can be rewritten as a sum of many small block matrices along the diagonal. This simple but important observation is very useful for our technical arguments. Define the block matrices

$$M_l^{*(m)} = (\sigma_{ij}^* I \{l \leq i < l + m, l \leq j < l + m\})_{p \times p}$$

and set

$$S^{*(m)} = \sum_{l=1-m}^p M_l^{*(m)}$$

for all integers  $1 - m \leq l \leq p$  and  $m \geq 1$ .

**Lemma 1** *The tapering estimator  $\hat{\Sigma}_k$  given in (4) can be written as*

$$\hat{\Sigma}_k = k_h^{-1} \left( S^{*(k)} - S^{*(k_h)} \right). \quad (6)$$

It is clear that the performance of the estimator  $\hat{\Sigma}_k$  depends on the choice of the tapering parameter  $k$ . The optimal choice of  $k$  critically depends on the norm under which the estimation error is measured. We will study in the next two sections the rate of convergence of the tapering estimator under both the operator norm and Frobenius norm. Together with the minimax lower bounds derived in Sections 3 and 4, the results show that a tapering estimator with the optimal choice of  $k$  attains the optimal rate of convergence under these two norms.

### 3 Rate Optimality under the Operator Norm

In this section we will establish the optimal rate of convergence under the operator norm. For  $1 \leq q \leq \infty$ , the matrix  $\ell_q$ -norm of a matrix  $A$  is defined by  $\|A\|_q = \max_{\|x\|_q=1} \|Ax\|_q$ . The commonly used operator norm  $\|\cdot\|$  coincides with the matrix  $\ell_2$ -norm  $\|\cdot\|_2$ . For a symmetric matrix  $A$ , it is known that the operator norm  $\|A\|$  is equal to the largest magnitude of eigenvalues of  $A$ . Hence it is also called the spectral norm. We will establish Theorem 1 by deriving a minimax upper bound using the tapering estimator and a matching minimax lower bound by a careful construction of a collection of multivariate normal distributions and the application of Assouad's Lemma and Le Cam's method. We shall focus on the case  $p \geq n^{\frac{1}{2\alpha+1}}$  in Sections 3.1 and 3.2. The case of  $p < n^{\frac{1}{2\alpha+1}}$ , which will be discussed in Section 3.3, is similar and slightly easier.

#### 3.1 Minimax Upper Bound under the Operator Norm

We derive in this section the risk upper bound for the tapering estimators defined in (6) under the operator norm. Throughout the paper we denote by  $C$  a generic positive constant which may vary from place to place but always depends only on indices  $\alpha$ ,  $M_0$  and  $M$  of the matrix family. We shall assume that the distribution of the  $X_i$ 's is subgaussian in the sense that there is  $\rho > 0$  such that

$$\mathbb{P}\{|\mathbf{v}^T(\mathbf{X}_1 - \mathbb{E}\mathbf{X}_1)| > t\} \leq e^{-t^2\rho/2} \text{ for all } t > 0 \text{ and } \|\mathbf{v}\|_2 = 1. \quad (7)$$

Let  $\mathcal{P}_\alpha = \mathcal{P}_\alpha(M_0, M, \rho)$  denote the set of distributions of  $\mathbf{X}_1$  that satisfy (3) and (7).

**Theorem 2** *The tapering estimator  $\hat{\Sigma}_k$ , defined in (6), of the covariance matrix  $\Sigma_{p \times p}$  with  $p \geq n^{\frac{1}{2\alpha+1}}$  satisfies*

$$\sup_{\mathcal{P}_\alpha} \mathbb{E} \left\| \hat{\Sigma}_k - \Sigma \right\|^2 \leq C \frac{k + \log p}{n} + Ck^{-2\alpha} \quad (8)$$

for  $k = o(n)$ ,  $\log p = o(n)$  and some constant  $C > 0$ . In particular, the estimator  $\hat{\Sigma} = \hat{\Sigma}_k$  with  $k = n^{\frac{1}{2\alpha+1}}$  satisfies

$$\sup_{\mathcal{P}_\alpha} \mathbb{E} \left\| \hat{\Sigma} - \Sigma \right\|^2 \leq Cn^{-\frac{2\alpha}{2\alpha+1}} + C \frac{\log p}{n}. \quad (9)$$

From (8) it is clear that the optimal choice of  $k$  is of order  $n^{\frac{1}{2\alpha+1}}$ . The upper bound given in (9) is thus rate optimal among the class of the tapering estimators defined in (6). The minimax lower bound derived in Section 3.2 shows that the estimator  $\hat{\Sigma}_k$  with  $k = n^{\frac{1}{2\alpha+1}}$  is in fact rate optimal among all estimators.

**Proof of Theorem 2:** Note that  $\Sigma^*$  is translation invariant and so is  $\hat{\Sigma}$ . We shall thus assume  $\mu = 0$  for the rest of the paper. Write

$$\Sigma^* = \frac{1}{n} \sum_{l=1}^n (\mathbf{X}_l - \bar{\mathbf{X}}) (\mathbf{X}_l - \bar{\mathbf{X}})^T = \frac{1}{n} \sum_{l=1}^n \mathbf{X}_l \mathbf{X}_l^T - \bar{\mathbf{X}} \bar{\mathbf{X}}^T$$

where  $\bar{\mathbf{X}} \bar{\mathbf{X}}^T$  is a higher order term (see Remark 1 at the end of this section). In what follows we shall ignore this negligible term and focus on the dominating term  $\frac{1}{n} \sum_{l=1}^n \mathbf{X}_l \mathbf{X}_l^T$ .

Set  $\tilde{\Sigma} = \frac{1}{n} \sum_{l=1}^n \mathbf{X}_l \mathbf{X}_l^T$  and write  $\tilde{\Sigma} = (\tilde{\sigma}_{ij})_{1 \leq i, j \leq p}$ . Let

$$\check{\Sigma} = (\check{\sigma}_{ij})_{1 \leq i, j \leq p} = (w_{ij} \tilde{\sigma}_{ij})_{1 \leq i, j \leq p} \quad (10)$$

with  $w_{ij}$  given in (5). Let  $\mathbf{X}_l = (X_1^l, X_2^l, \dots, X_p^l)^T$ . We then write  $\tilde{\sigma}_{ij} = \frac{1}{n} \sum_{l=1}^n X_i^l X_j^l$ . It is easy to see

$$\mathbb{E} \tilde{\sigma}_{ij} = \sigma_{ij} \quad (11)$$

$$\text{Var}(\tilde{\sigma}_{ij}) = \frac{1}{n} \text{Var}(X_i^l X_j^l) \leq \frac{1}{n} \mathbb{E}(X_i^l X_j^l)^2 \leq \frac{1}{n} \mathbb{E}(X_i^l)^2 \mathbb{E}(X_j^l)^2 \leq \frac{C}{n} \quad (12)$$

i.e.,  $\tilde{\sigma}_{ij}$  is an unbiased estimator of  $\sigma_{ij}$  with a variance  $O(1/n)$ .

We will first show that the variance part satisfies

$$\mathbb{E} \left\| \check{\Sigma} - \mathbb{E} \check{\Sigma} \right\|^2 \leq C \frac{k + \log p}{n} \quad (13)$$

and the bias part satisfies

$$\left\| \mathbb{E} \check{\Sigma} - \Sigma \right\|^2 \leq Ck^{-2\alpha}. \quad (14)$$

It then follows immediately that

$$\mathbb{E} \left\| \check{\Sigma} - \Sigma \right\|^2 \leq 2\mathbb{E} \left\| \check{\Sigma} - \mathbb{E}\check{\Sigma} \right\|^2 + 2 \left\| \mathbb{E}\check{\Sigma} - \Sigma \right\|^2 \leq 2C \left( \frac{k + \log p}{n} + k^{-2\alpha} \right).$$

This proves (8) and equation (9) then follows. Since  $p \geq n^{\frac{1}{2\alpha+1}}$ , we may choose

$$k = n^{\frac{1}{2\alpha+1}} \tag{15}$$

and the estimator  $\hat{\Sigma}$  with  $k$  given in (15) satisfies

$$\mathbb{E} \left\| \hat{\Sigma} - \Sigma \right\|^2 \leq 2C \left( n^{-\frac{2\alpha}{2\alpha+1}} + \frac{\log p}{n} \right).$$

Theorem 2 is then proved.

We first prove the risk upper bound (14) for the bias part. It is well known that the operator norm of a symmetric matrix  $A = (a_{ij})_{p \times p}$  is bounded by its  $\ell_1$  norm, i.e.,

$$\|A\| \leq \|A\|_1 = \max_{i=1, \dots, p} \sum_{j=1}^p |a_{ij}|.$$

((see e.g. page 15 in Golub and Van Loan (1983)). This result was used in Bickel and Levina (2008a, b) to obtain rates of convergence for their proposed procedures under the operator norm (see discussions in Section 3.3). We bound the operator norm of the bias part  $\mathbb{E}\check{\Sigma} - \Sigma$  by its  $\ell_1$  norm. Since  $\mathbb{E}\check{\sigma}_{ij} = \sigma_{ij}$ , we have

$$\mathbb{E}\check{\Sigma} - \Sigma = ((w_{ij} - 1) \sigma_{ij})_{p \times p}$$

where  $w_{ij} \in [0, 1]$  and is exactly 1 when  $|i - j| \leq k$ , then

$$\left\| \mathbb{E}\check{\Sigma} - \Sigma \right\|^2 \leq \left[ \max_{i=1, \dots, p} \sum_{j: |i-j| > k} |\sigma_{ij}| \right]^2 \leq M^2 k^{-2\alpha}.$$

Now we establish (13) which is relatively complicated. The key idea in the proof is to write the whole matrix as an average of matrices which are sum of a large number of small disjoint block matrices, and for each small block matrix the classical random matrix theory can be applied. The following lemma shows that the operator norm of the random matrix  $\check{\Sigma} - \mathbb{E}\check{\Sigma}$  is controlled by the maximum of operator norms of  $p$  number of  $k \times k$  random matrices. Let  $M_l^{(m)} = (\tilde{\sigma}_{ij} I \{l \leq i < l + m, l \leq j < l + m\})_{p \times p}$ . Define

$$N_l^{(m)} = \max_{1 \leq l \leq p-m+1} \left\| M_l^{(m)} - \mathbb{E}M_l^{(m)} \right\|.$$

**Lemma 2** Let  $\check{\Sigma}$  be defined as in (6). Then

$$\left\| \check{\Sigma} - \mathbb{E}\check{\Sigma} \right\| \leq 3N_l^{(m)}.$$

For each small  $m \times m$  random matrix with  $m = k$ , we control its operator norm as follows.

**Lemma 3** There is a constant  $\rho_1 > 0$  such that

$$\mathbb{P} \left\{ N_l^{(m)} > x \right\} \leq 2p5^m \exp(-nx^2\rho_1) \quad (16)$$

for all  $0 < x < \rho_1$  and  $1 - m \leq l \leq p$ .

With Lemmas 2 and 3 we are now ready to show the variance bound (13). By Lemma 2 we have

$$\begin{aligned} \mathbb{E} \left\| \check{\Sigma} - \mathbb{E}\check{\Sigma} \right\|^2 &\leq 9\mathbb{E} \left( N_l^{(m)} \right)^2 = 9\mathbb{E} \left( N_l^{(m)} \right)^2 \left[ I \left( N_l^{(m)} \leq x \right) + I \left( N_l^{(m)} > x \right) \right] \\ &\leq 9 \left[ x^2 + \mathbb{E} \left( N_l^{(m)} \right)^2 I \left( N_l^{(m)} > x \right) \right]. \end{aligned}$$

Note that  $\left\| \mathbb{E}\check{\Sigma} \right\| \leq \|\Sigma\|$ , which is bounded by a constant, and  $\left\| \check{\Sigma} \right\| \leq \left\| \check{\Sigma} \right\|_F$ . The Cauchy–Schwarz inequality then implies

$$\begin{aligned} \mathbb{E} \left\| \check{\Sigma} - \mathbb{E}\check{\Sigma} \right\|^2 &\leq C_1 \left[ x^2 + \mathbb{E} \left( \left\| \check{\Sigma} \right\|_F^2 + C \right) I \left( N_l^{(m)} > x \right) \right] \\ &\leq C_1 \left[ x^2 + \sqrt{\mathbb{E} \left( \left\| \check{\Sigma} \right\|_F + C \right)^4} \sqrt{\mathbb{P} \left( N_l^{(m)} > x \right)} \right]. \end{aligned}$$

Set  $x = 4\sqrt{\frac{\log p + m}{n\rho_1}}$ . Then  $x$  is bounded by  $\rho_1$  as  $n \rightarrow \infty$ . From Lemma 3 we obtain

$$\mathbb{E} \left\| \check{\Sigma} - \mathbb{E}\check{\Sigma} \right\|^2 \leq C \left[ \frac{\log p + m}{n} + p^2 \cdot (p5^m \cdot p^{-8}e^{-8m})^{1/2} \right] \leq C_1 \left( \frac{\log p + m}{n} \right). \quad \blacksquare \quad (17)$$

**Remark 1** In the proof of Theorem 2, the term  $\bar{\mathbf{X}}\bar{\mathbf{X}}^T$  was ignored. It is not difficult to see that this term has negligible contribution after tapering. Let  $H = \bar{\mathbf{X}}\bar{\mathbf{X}}^T$  and  $H = (h_{ij})_{p \times p}$ . Define  $H_l^{(m)} = (h_{ij}I \{l \leq i < l + m, l \leq j < l + m\})_{p \times p}$ . Similarly to Lemma 3 it can be shown that

$$\mathbb{P} \left\{ \max_{1 \leq l \leq p-m+1} \left\| H_l^{(m)} - \mathbb{E}H_l^{(m)} \right\| > t \right\} \leq 2p5^m \exp(-nt\rho_2) \quad (18)$$

for all  $0 < t < \rho_2$  and  $1 - m \leq l \leq p$ . Note that  $\mathbb{E}H = \frac{1}{n}\Sigma$ , then

$$\mathbb{E} \|H\|^2 \leq 2\mathbb{E} \|H - \mathbb{E}H\|^2 + 2\|\mathbb{E}H\|^2 \leq 2\mathbb{E} \|H - \mathbb{E}H\|^2 + 2M_0^2/n^2.$$

Let  $t = 16 \frac{\log p + m}{n \rho^2}$ . From equation (18) we have

$$\begin{aligned} \mathbb{E} \|H - \mathbb{E}H\|^2 &\leq t^2 + \mathbb{E} \|H - \mathbb{E}H\|^2 I \left( \max_{1 \leq l \leq p-m+1} \|H_l^{(m)} - \mathbb{E}H_l^{(m)}\| > t \right) \\ &= t^2 + o(t^2) \leq C \left( \frac{\log p + m}{n} \right)^2 \end{aligned}$$

by similar arguments as for equation (17). Therefore  $H$  has a negligible contribution to the risk.

### 3.2 Lower Bound under the Operator Norm

Theorem 2 in Section 3.1 shows that the optimal tapering estimator attains the rate of convergence  $n^{-\frac{2\alpha}{2\alpha+1}} + \frac{\log p}{n}$ . In this section we shall show that this rate of convergence is indeed optimal among all estimators by showing that the upper bound in equation (9) can not be improved. More specifically we shall show that the following minimax lower bound holds.

**Theorem 3** *Suppose  $p \leq \exp(\gamma n)$  for some constant  $\gamma > 0$ . The minimax risk for estimating the covariance matrix  $\Sigma$  over  $\mathcal{P}_\alpha$  under the operator norm satisfies*

$$\inf_{\hat{\Sigma}} \sup_{\mathcal{P}_\alpha} \mathbb{E} \left\| \hat{\Sigma} - \Sigma \right\|^2 \geq cn^{-\frac{2\alpha}{2\alpha+1}} + c \frac{\log p}{n}.$$

The basic strategy underlying the proof of Theorem 3 is to carefully construct a finite collection of multivariate normal distributions and calculate the total variation affinity between pairs of probability measures in the collection.

We shall now define a parameter space that is appropriate for the minimax lower bound argument. For given positive integers  $k$  and  $m$  with  $2k \leq p$  and  $1 \leq m \leq k$ , define the  $p \times p$  matrix  $B(m, k) = (b_{ij})_{p \times p}$  with

$$b_{ij} = I \{i = m \text{ and } m + 1 \leq j \leq 2k, \text{ or } j = m \text{ and } m + 1 \leq i \leq 2k\}.$$

Set  $k = n^{\frac{1}{2\alpha+1}}$  and  $a = k^{-(\alpha+1)}$ . We then define the collection of  $2^k$  covariance matrices as

$$\mathcal{F}_{11} = \left\{ \Sigma(\theta) : \Sigma(\theta) = I_p + \tau a \sum_{m=1}^k \theta_m B(m, k), \quad \theta = (\theta_m) \in \{0, 1\}^k \right\} \quad (19)$$

where  $I_p$  is the  $p \times p$  identity matrix and  $0 < \tau < 2^{-\alpha-1}M$ . Without loss of generality we assume that  $M_0 > 1$  and  $\rho > 1$ . Otherwise we replace  $I_p$  in equation (19) by  $\varepsilon I_p$  for

$0 < \varepsilon < \min \{M_0, \rho\}$ . For  $0 < \tau < 2^{-\alpha-1}M$  it is easy to check that  $\mathcal{F}_{11} \subset \mathcal{F}_\alpha(M_0, M)$  as  $n \rightarrow \infty$ . In addition to  $\mathcal{F}_{11}$  we also define a collection of diagonal matrices

$$\mathcal{F}_{12} = \left\{ \Sigma_m : \Sigma_m = I_p + \left( \sqrt{\frac{\tau}{n} \log p_1} I \{i = j = m\} \right)_{p \times p}, 0 \leq m \leq p_1 \right\} \quad (20)$$

where  $p_1 = \min \{p, e^{n/2}\}$  and  $0 < \tau < \min \{(M_0 - 1)^2, (\rho - 1)^2, 1\}$ . Let  $\mathcal{F}_1 = \mathcal{F}_{11} \cup \mathcal{F}_{12}$ . It is clear that  $\mathcal{F}_1 \subset \mathcal{F}_\alpha(M_0, M)$ .

We shall show below separately that the minimax risks over multivariate normal distributions with covariance matrix in (19) and (20) satisfy

$$\inf_{\hat{\Sigma}} \sup_{\mathcal{F}_{11}} \mathbb{E} \left\| \hat{\Sigma} - \Sigma \right\|^2 \geq cn^{-\frac{2\alpha}{2\alpha+1}} \quad (21)$$

and

$$\inf_{\hat{\Sigma}} \sup_{\mathcal{F}_{12}} \mathbb{E} \left\| \hat{\Sigma} - \Sigma \right\|^2 \geq c \frac{\log p}{n} \quad (22)$$

for some constant  $c > 0$ . Equations (21) and (22) together imply

$$\inf_{\hat{\Sigma}} \sup_{\mathcal{F}_1} \mathbb{E} \left\| \hat{\Sigma} - \Sigma \right\|^2 \geq \frac{c}{2} \left( n^{-\frac{2\alpha}{2\alpha+1}} + \frac{\log p}{n} \right) \quad (23)$$

for multivariate normal distributions and this proves Theorem 3. We shall establish the lower bound (21) by using Assouad's Lemma in Section 3.2.1 and the lower bound (22) by using Le Cam's method and a two-point argument in Section 3.2.2.

### 3.2.1 A Lower Bound by Assouad's Lemma

The key technical tool to establish equation (21) is the Assouad's lemma in Assouad (1983). It gives a lower bound for the maximum risk over the parameter set  $\Theta = \{0, 1\}^k$  to the problem of estimating an arbitrary quantity  $\psi(\theta)$ , belonging to a metric space with metric  $d$ . Let  $H(\theta, \theta') = \sum_{i=1}^k |\theta_i - \theta'_i|$  be the Hamming distance on  $\{0, 1\}^k$ , which counts the number of positions at which  $\theta$  and  $\theta'$  differ. For two probability measures  $P$  and  $Q$  with density  $p$  and  $q$  with respect to any common dominating measure  $\mu$ , write the total variation affinity  $\|P \wedge Q\| = \int p \wedge q d\mu$ . Assouad's Lemma provides a minimax lower bound for estimating  $\psi(\theta)$ .

**Lemma 4 (Assouad)** *Let  $\Theta = \{0, 1\}^k$  and let  $T$  be an estimator based on an observation from a distribution in the collection  $\{P_\theta, \theta \in \Theta\}$ . Then for all  $s > 0$*

$$\max_{\theta \in \Theta} 2^s \mathbb{E}_\theta d^s(T, \psi(\theta)) \geq \min_{H(\theta, \theta') \geq 1} \frac{d^s(\psi(\theta), \psi(\theta'))}{H(\theta, \theta')} \cdot \frac{k}{2} \cdot \min_{H(\theta, \theta')=1} \|\mathbb{P}_\theta \wedge \mathbb{P}_{\theta'}\|.$$

The Assouad's lemma is connected to multiple comparisons. In total there are  $k$  comparisons. The lower bound has three factors. The first factor is basically the minimum cost of making a mistake per comparison, and the last factor is the lower bound for the total probability of making type I and type II errors for each comparison, and  $k/2$  is the expected number of mistakes one makes when  $\mathbb{P}_\theta$  and  $\mathbb{P}_{\theta'}$  are not distinguishable from each other when  $H(\theta, \theta') = 1$ .

We now prove the lower bound (21). Let  $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{iid}{\sim} N(0, \Sigma(\theta))$  with  $\Sigma(\theta) \in \mathcal{F}_{11}$ . Denote the joint distribution by  $P_\theta$ . Applying Assouad's Lemma to the parameter space  $\mathcal{F}_{11}$ , we have

$$\inf_{\hat{\Sigma}} \max_{\theta \in \{0,1\}^k} 2^2 E_\theta \left\| \hat{\Sigma} - \Sigma(\theta) \right\|^2 \geq \min_{H(\theta, \theta') \geq 1} \frac{\|\Sigma(\theta) - \Sigma(\theta')\|^2 k}{H(\theta, \theta')} \frac{k}{2} \min_{H(\theta, \theta')=1} \|P_\theta \wedge P_{\theta'}\|. \quad (24)$$

We shall state the bounds for the the first and third factors on the right hand of (24) in two lemmas. The proof of these lemmas is given in Section 7.

**Lemma 5** *Let  $\Sigma(\theta)$  be defined as in (19). Then for some constant  $c > 0$*

$$\min_{H(\theta, \theta') \geq 1} \frac{\|\Sigma(\theta) - \Sigma(\theta')\|^2}{H(\theta, \theta')} \geq cka^2.$$

**Lemma 6** *Let  $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{iid}{\sim} N(0, \Sigma(\theta))$  with  $\Sigma(\theta) \in \mathcal{F}_{11}$ . Denote the joint distribution by  $P_\theta$ . Then for some constant  $c > 0$*

$$\min_{H(\theta, \theta')=1} \|P_\theta \wedge P_{\theta'}\| \geq c.$$

It then follows from Lemmas 5 and 6 together with the fact  $k = n^{\frac{1}{2\alpha+1}}$

$$\max_{\Sigma(\theta) \in \mathcal{F}_{11}} 2^2 E_\theta \left\| \hat{\Sigma} - \Sigma(\theta) \right\|^2 \geq \frac{c^2}{2} k^2 a^2 \geq c_1 n^{-\frac{2\alpha}{2\alpha+1}}$$

for some  $c_1 > 0$ . ■

### 3.2.2 A Lower Bound Using Le Cam's Method

We now apply Le Cam's method to derive the lower bound (22) for the minimax risk. Let  $X$  be an observation from a distribution in the collection  $\{P_\theta, \theta \in \Theta\}$  where  $\Theta = \{\theta_0, \theta_1, \dots, \theta_{p_1}\}$ . Le Cam's method, which is based on a two-point testing argument, gives a lower bound for the maximum estimation risk over the parameter set  $\Theta$ . More specifically, let  $L$  be the loss function. Define  $r(\theta_0, \theta_m) = \inf_t [L(t, \theta_0) + L(t, \theta_m)]$  and  $r_{\min} = \inf_{1 \leq m \leq p_1} r(\theta_0, \theta_m)$ , and denote  $\bar{\mathbb{P}} = \frac{1}{p_1} \sum_{m=1}^{p_1} \mathbb{P}_{\theta_m}$ .

**Lemma 7** *Let  $T$  be an estimator of  $\theta$  based on an observation from a distribution in the collection  $\{P_\theta, \theta \in \Theta = \{\theta_0, \theta_1, \dots, \theta_{p_1}\}\}$ , then*

$$\sup_{\theta} \mathbb{E}L(T, \theta) \geq \frac{1}{2} r_{\min} \|\mathbb{P}_{\theta_0} \wedge \bar{\mathbb{P}}\|$$

We refer to Yu (1997) for more detailed discussions on Le Cam's method.

To apply Le Cam's method, we need to first construct a parameter set. For  $1 \leq m \leq p_1$ , let  $\Sigma_m$  be a diagonal covariance matrix with  $\sigma_{mm} = 1 + \sqrt{\tau \frac{\log p_1}{n}}$ ,  $\sigma_{ii} = 1$  for  $i \neq m$ , and let  $\Sigma_0$  be the identity matrix. Let  $\mathbf{X}_l = (X_1^l, X_2^l, \dots, X_{p_1}^l)^T \sim N(0, \Sigma_m)$ , and denote the joint density of  $\mathbf{X}_1, \dots, \mathbf{X}_n$  by  $f_m$ ,  $0 \leq m \leq p_1$  with  $p_1 = \max\{p, e^{n/2}\}$ , which can be written as follows

$$f_m = \prod_{1 \leq i \leq n, 1 \leq j \leq p, j \neq m} \phi_1(x_j^i) \cdot \prod_{1 \leq i \leq n} \phi_{\sigma_{mm}}(x_m^i)$$

where  $\phi_\sigma$ ,  $\sigma = 1$  or  $\sigma_{mm}$ , is the density of  $N(0, \sigma^2)$ . Denote by  $f_0$  the joint density of  $\mathbf{X}_1, \dots, \mathbf{X}_n$  when  $\mathbf{X}_l \sim N(0, \Sigma_0)$ .

Let  $\theta_m = \Sigma_m$  for  $0 \leq m \leq p_1$  and the loss function  $L$  be the squared operator norm. It is easy to see  $r(\theta_0, \theta_m) = \frac{1}{2} \tau \frac{\log p_1}{n}$  for all  $1 \leq m \leq p_1$ . Then the lower bound (22) follows immediately from Lemma 7 if there is a constant  $c > 0$  such that

$$\|\mathbb{P}_{\theta_0} \wedge \bar{\mathbb{P}}\| \geq c. \quad (25)$$

Note that for any two densities  $q_0$  and  $q_1$ ,  $\int q_0 \wedge q_1 d\mu = 1 - \frac{1}{2} \int |q_0 - q_1| d\mu$ , and the Jensen's inequality implies

$$\left[ \int |q_0 - q_1| d\mu \right]^2 = \left( \int \left| \frac{q_0 - q_1}{q_1} \right| q_1 d\mu \right)^2 \leq \int \frac{(q_0 - q_1)^2}{q_1} d\mu = \int \frac{q_0^2}{q_1} d\mu - 1.$$

Hence  $\int q_0 \wedge q_1 d\mu \geq 1 - \frac{1}{2} \left( \int \frac{q_0^2}{q_1} d\mu - 1 \right)^{1/2}$ . To establish equation (25), it thus suffices to show that  $\int \left( \frac{1}{p_1} \sum_{m=1}^{p_1} f_m \right)^2 / f_0 d\mu - 1 \rightarrow 0$ , i.e.,

$$\frac{1}{p_1^2} \sum_{m=1}^{p_1} \int \frac{f_m^2}{f_0} d\mu + \frac{1}{p_1^2} \sum_{m \neq j} \int \frac{f_m f_j}{f_0} d\mu - 1 \rightarrow 0. \quad (26)$$

We now calculate  $\int \frac{f_m f_j}{f_0} d\mu$ . For  $m \neq j$  it is easy to see

$$\int \frac{f_m f_j}{f_0} d\mu - 1 = 0.$$

When  $m = j$ , we have

$$\begin{aligned} \int \frac{f_m^2}{f_0} d\mu &= \frac{(\sqrt{2\pi}\sigma_{mm})^{-2n}}{(\sqrt{2\pi})^{-n}} \prod_{1 \leq i \leq n} \int \exp \left[ (x_m^i)^2 \left( -\frac{1}{\sigma_{mm}} + \frac{1}{2} \right) \right] dx_m^i \\ &= \left[ 1 - (1 - \sigma_{mm})^2 \right]^{-n/2} = \left( 1 - \tau \frac{\log p_1}{n} \right)^{-n/2}. \end{aligned}$$

Thus

$$\begin{aligned} \int \left( \frac{1}{p_1} \sum_{m=1}^{p_1} f_m \right)^2 / f_0 d\mu - 1 &= \frac{1}{p_1^2} \sum_{m=1}^{p_1} \left( \int \frac{f_m^2}{f_0} d\mu - 1 \right) \leq \frac{1}{p_1} \left( 1 - \tau \frac{\log p_1}{n} \right)^{-n/2} - \frac{1}{p_1} \\ &= \exp \left[ -\log p_1 - \frac{n}{2} \log \left( 1 - \tau \frac{\log p_1}{n} \right) \right] - \frac{1}{p_1} \rightarrow 0 \quad (27) \end{aligned}$$

for  $0 < \tau < 1$ , where the last step follows from the inequality  $\log(1-x) \geq -2x$  for  $0 < x < 1/2$ . Equation (27) together with Lemma 7 now immediately imply the lower bound given in equation (22). ■

**Remark 2** In covariance matrix estimation literature, it is commonly assumed that  $\frac{\log p}{n} \rightarrow 0$ . See, for example, Bickel and Levina (2008a). The lower bound given in this section implies that this assumption is necessary for estimating the covariance matrix consistently under the operator norm.

### 3.3 Discussion

Theorems 2 and 3 together shows that the minimax risk for estimating the covariance matrices over the distribution space  $\mathcal{P}_\alpha$  satisfies, for  $p \geq n^{\frac{1}{2\alpha+1}}$ ,

$$\inf_{\hat{\Sigma}} \sup_{\mathcal{P}_\alpha} \mathbb{E} \left\| \hat{\Sigma} - \Sigma \right\|^2 \asymp n^{-\frac{2\alpha}{2\alpha+1}} + \frac{\log p}{n}. \quad (28)$$

The results also show that the tapering estimator  $\hat{\Sigma}_k$  with tapering parameter  $k = n^{\frac{1}{2\alpha+1}}$  attains the optimal rate of convergence  $n^{-\frac{2\alpha}{2\alpha+1}} + \frac{\log p}{n}$ .

A few interesting points can be made on the optimal rate of convergence  $n^{-\frac{2\alpha}{2\alpha+1}} + \frac{\log p}{n}$ . When the dimension  $p$  is relatively small, i.e.,  $\log p = o(n^{\frac{1}{2\alpha+1}})$ ,  $p$  has no effect on the convergence rate and the rate is purely driven by the ‘‘smoothness’’ parameter  $\alpha$ . However, when  $p$  is large, i.e.,  $\log p \gg n^{\frac{1}{2\alpha+1}}$ ,  $p$  plays a significant role in determining the minimax rate.

We should emphasize that the optimal choice of the tapering parameter  $k \asymp n^{\frac{1}{2\alpha+1}}$  is different from the optimal choice for estimating the rows/columns as vectors under mean

squared error loss. Straightforward calculation shows that in the latter case the best cutoff is  $k \asymp n^{\frac{1}{2(\alpha+1)}}$  so that the tradeoff between the squared bias and the variance is optimal. With  $k \asymp n^{\frac{1}{2\alpha+1}}$ , the tapering estimator has smaller squared bias than the variance as a vector estimator of each row/column.

It is also interesting to compare our results with those given in Bickel and Levina (2008a). A banding estimator with bandwidth  $k = \left(\frac{\log p}{n}\right)^{\frac{1}{2(\alpha+1)}}$  was proposed and the rate of convergence  $\left(\frac{\log p}{n}\right)^{\frac{\alpha}{\alpha+1}}$  was proved. It is easy to see that the banding estimator given in Bickel and Levina (2008a) is not rate optimal. Take for example  $\alpha = 1/2$  and  $p = e^{\sqrt{n}}$ . Their rate is  $n^{-\frac{1}{6}}$ , while the optimal rate in Theorem 1 is  $n^{-\frac{1}{2}}$ .

It is instructive to take a closer look at the motivation behind the construction of the banding estimator in Bickel and Levina (2008a). Let the banding estimator be

$$\hat{\Sigma}_B = (\sigma_{ij}^* I \{|i - j| \leq k\}) \quad (29)$$

and denote  $\hat{\Sigma}_B - \mathbb{E}\hat{\Sigma}_B$  by  $V$ , and let  $V = (v_{ij})$ . An important step in the proof of Theorem 1 in Bickel and Levina (2008a) is to control the operator norm by the  $\ell_1$  norm as follows,

$$\begin{aligned} \mathbb{E} \left\| \hat{\Sigma}_B - \mathbb{E}\hat{\Sigma}_B \right\|^2 &\leq \mathbb{E} \left\| \hat{\Sigma}_B - \mathbb{E}\hat{\Sigma}_B \right\|_1^2 = \mathbb{E} \left( \max_{j=1, \dots, p} \sum_i |v_{ij}| \right)^2 \\ &\leq C \left( \frac{k}{\sqrt{n}} \sqrt{\log p} \right)^2 = C \frac{k^2 \log p}{n}. \end{aligned}$$

Note that  $\mathbb{E}[|v_{ij}| I \{|i - j| \leq k\}] \asymp 1/\sqrt{n}$ , then  $\mathbb{E} \sum_i |v_{ij}| \asymp k/\sqrt{n}$ . It is then expected that  $\mathbb{E} (\max_{j=1, \dots, p} \sum_i |v_{ij}|)^2 \leq C \left( \frac{k}{\sqrt{n}} \sqrt{\log p} \right)^2$  (see Bickel and Levina (2008a) for details) and so

$$\mathbb{E} \left\| \check{\Sigma} - \Sigma \right\|_1^2 \leq C \frac{k^2 \log p}{n} + Ck^{-2\alpha}$$

An optimal tradeoff of  $k$  is then  $\left(\frac{\log p}{n}\right)^{\frac{1}{2(\alpha+1)}}$  which implies a rate of  $\left(\frac{\log p}{n}\right)^{-\frac{\alpha}{\alpha+1}}$  in Theorem 1 in Bickel and Levina (2008a). This rate is slower than the optimal rate  $n^{-\frac{2\alpha}{2\alpha+1} + \frac{\log p}{n}}$  in Theorem 1.

We have considered the parameter space  $\mathcal{F}_\alpha$  defined in (3). Other similar parameter spaces can also be considered. For example in time series analysis it is often assumed the covariance  $|\sigma_{ij}|$  decays at the rate  $|i - j|^{-(\alpha+1)}$  for some  $\alpha > 0$ . Consider the collection of positive definite symmetric matrices satisfying the following conditions:

$$\mathcal{G}_\alpha = \mathcal{G}_\alpha(M_0, M_1) = \left\{ \Sigma : |\sigma_{ij}| \leq M_1 |i - j|^{-(\alpha+1)} \text{ for } i \neq j \text{ and } \lambda_{\max}(\Sigma) \leq M_0 \right\} \quad (30)$$

where  $\lambda_{\max}(\Sigma)$  is the maximum eigenvalues of the matrix  $\Sigma$ . Note that  $\mathcal{G}_\alpha(M_0, M_1)$  is a subset of  $\mathcal{F}_\alpha(M_0, M)$  as long as  $M_1 \leq \alpha M$ . Using virtually identical arguments one can show that

$$\inf_{\hat{\Sigma}} \sup_{\mathcal{P}'_\alpha} \mathbb{E} \left\| \hat{\Sigma} - \Sigma \right\|^2 \asymp n^{-\frac{2\alpha}{2\alpha+1}} + \frac{\log p}{n}$$

Let  $\mathcal{P}'_\alpha = \mathcal{P}'_\alpha(M_0, M, \rho)$  denote the set of distributions of  $\mathbf{X}_1$  that satisfies (7) and (30).

**Remark 3** Both the tapering estimator proposed in this paper and banding estimator given in Bickel and Levina (2008a) are not necessarily positive-semidefinite. A practical proposal to avoid this would be to project the estimator  $\hat{\Sigma}$  to the space of positive-semidefinite matrices under the operator norm. More specifically, one may first diagonalize  $\hat{\Sigma}$  and then replace negative eigenvalues by 0. The resulting estimator is then positive-semidefinite.

### 3.3.1 The Case of $p < n^{\frac{1}{2\alpha+1}}$

We have focused on the case  $p \geq n^{\frac{1}{2\alpha+1}}$  in Sections 3.1 and 3.2. The case of  $p < n^{\frac{1}{2\alpha+1}}$  can be handled in a similar way. The main difference is that in this case we no longer have a tapering estimator  $\hat{\Sigma}_k$  with  $k = n^{\frac{1}{2\alpha+1}}$  because  $k > p$ . Instead the maximum likelihood estimator  $\Sigma^*$  can be used directly. It is easy to show in this case,

$$\sup_{\mathcal{P}_\alpha} \mathbb{E} \|\Sigma^* - \Sigma\|^2 \leq C \frac{p}{n}. \quad (31)$$

The lower bound can also be obtained by the application of Assouad's Lemma and by using a parameter space that is similar to  $\mathcal{F}_{11}$ . To be more specific, For an integer  $1 \leq m \leq p/2$ , define the  $p \times p$  matrix  $B_m = (b_{ij})_{p \times p}$  with

$$b_{ij} = I \{i = m \text{ and } m + 1 \leq j \leq p, \text{ or } j = m \text{ and } m + 1 \leq i \leq p\}.$$

Define the collection of  $2^{p/2}$  covariance matrices as

$$\mathcal{F}^* = \left\{ \Sigma(\theta) : \Sigma(\theta) = I_p + \tau \frac{1}{\sqrt{np}} \sum_{m=1}^{p/2} \theta_m B(m, k), \quad \theta = (\theta_m) \in \{0, 1\}^{p/2} \right\}. \quad (32)$$

Since  $p < n^{\frac{1}{2\alpha+1}}$ , then  $\frac{1}{\sqrt{np}} < 2^{\alpha+1/2} p^{-(\alpha+1)}$ . Again it is easy to check  $\mathcal{F}^* \subset \mathcal{F}_\alpha(M_0, M)$  when  $0 < \tau < 2^{-\alpha-1} M$ . The following lower bound then follows from the same argument as in Section 3.2.1

$$\inf_{\hat{\Sigma}} \sup_{\mathcal{F}^*} \mathbb{E} \left\| \hat{\Sigma} - \Sigma \right\|^2 \geq cp \left( \frac{1}{\sqrt{np}} \right)^2 \cdot \frac{p}{2} \cdot c_1 \geq c_2 \frac{p}{n}. \quad (33)$$

Equations (31) and (33) together yield the minimax rate of convergence for the case  $p \leq n^{\frac{1}{2\alpha+1}}$ ,

$$\inf_{\hat{\Sigma}} \sup_{\mathcal{P}_\alpha} \mathbb{E} \left\| \hat{\Sigma} - \Sigma \right\|^2 \asymp \frac{p}{n}. \quad (34)$$

This together with equation (28) give the optimal rate of convergence

$$\inf_{\hat{\Sigma}} \sup_{\mathcal{P}_\alpha} \mathbb{E} \left\| \hat{\Sigma} - \Sigma \right\|^2 \asymp \min \left\{ n^{-\frac{2\alpha}{2\alpha+1}} + \frac{\log p}{n}, \frac{p}{n} \right\}. \quad (35)$$

## 4 Rate Optimality under the Frobenius Norm

In addition to the operator norm, the Frobenius norm is another commonly used matrix norm. The Frobenius norm is used in defining the numerical rank of a matrix which is useful in many applications, such as the principle component analysis. See, for e.g., Rudelson and Vershynin (2007). The Frobenius norm has also been used in the literature for measuring the accuracy of a covariance matrix estimator. See, e.g., Lam and Fan (2007) and Ravikumar, et al. (2008). In this section we consider the optimal rate of convergence for covariance matrix estimation under the Frobenius norm. The Frobenius norm of a matrix  $A = (a_{ij})$  is defined as the  $\ell_2$  vector norm of all entries in the matrix,

$$\|A\|_F = \sqrt{\sum_{i,j} a_{ij}^2}.$$

This is equivalent to treating the matrix  $A$  as a vector of length  $p^2$ . It is easy to see that the operator norm is bounded by the Frobenius norm, i.e.,  $\|A\| \leq \|A\|_F$ .

The following theorem gives the minimax rate of convergence for estimating the covariance matrix  $\Sigma$  under the Frobenius norm based on the sample  $\{\mathbf{X}_1, \dots, \mathbf{X}_n\}$ .

**Theorem 4** *The minimax risk under the Frobenius norm satisfies*

$$\inf_{\hat{\Sigma}} \sup_{\mathcal{P}_\alpha} \mathbb{E} \frac{1}{p} \left\| \hat{\Sigma} - \Sigma \right\|_F^2 \asymp \min \left\{ n^{-\frac{2\alpha+1}{2(\alpha+1)}}, \frac{p}{n} \right\} \quad (36)$$

and

$$\inf_{\hat{\Sigma}} \sup_{\mathcal{P}_\alpha} \mathbb{E} \frac{1}{p} \left\| \hat{\Sigma} - \Sigma \right\|_F^2 \asymp \min \left\{ n^{-\frac{2\alpha}{2\alpha+1}}, \frac{p}{n} \right\}$$

We shall establish below separately the minimax upper bound and minimax lower bound.

## 4.1 Upper Bound under the Frobenius Norm

We will only prove the upper bound for the distribution set  $\mathcal{P}'_\alpha$  given in (30). The proof for the parameter space  $\mathcal{P}_\alpha$  is similar. The minimax upper bound is derived by again considering the tapering estimator (4). Under the Frobenius norm the risk function is separable. The risk of the tapering estimator can be bounded separately under the squared  $\ell_2$  loss for each row/column. This method has been commonly used in nonparametric function estimation using orthogonal basis expansions. Since

$$\mathbb{E}\tilde{\sigma}_{ij} = \sigma_{ij}, \text{ and } \text{Var}(\tilde{\sigma}_{ij}) \leq \frac{C}{n}$$

for the tapering estimator (4), we have

$$\mathbb{E} (w_{ij}\tilde{\sigma}_{ij} - \sigma_{ij})^2 \leq (1 - w_{ij})^2 \sigma_{ij}^2 + w_{ij}^2 \frac{C}{n}.$$

It can be seen easily that

$$\frac{1}{p} \mathbb{E} \left\| \check{\Sigma} - \Sigma \right\|_F^2 \leq \frac{1}{p} \sum_{\{(i,j):k_h < |i-j|\}} \sigma_{ij}^2 + \frac{1}{p} \sum_{\{(i,j):|i-j| \leq k\}} \left[ (1 - w_{ij})^2 \sigma_{ij}^2 + w_{ij}^2 \frac{C}{n} \right] \equiv R_1 + R_2$$

The assumption  $\lambda_{\max}(\Sigma) \leq M_0$  implies that  $\sigma_{ii} \leq M_0$  for all  $i$ . Since  $|\sigma_{ij}|$  is also uniformly bounded for all  $i \neq j$  from Assumption (30), we immediately have  $R_2 \leq C \frac{k}{n}$ .

It is easy to show that

$$\frac{1}{p} \sum_{\{(i,j):k < |i-j|\}} \sigma_{ij}^2 \leq C k^{-2\alpha-1} \quad (37)$$

where  $|\sigma_{ij}| \leq C_1 |i-j|^{-(\alpha+1)}$  for all  $i \neq j$ . Thus

$$\mathbb{E} \frac{1}{p} \left\| \check{\Sigma} - \Sigma \right\|_F^2 \leq C k^{-2\alpha-1} + C \frac{k}{n} \leq C_2 n^{-\frac{2\alpha+1}{2(\alpha+1)}} \quad (38)$$

by choosing

$$k = n^{\frac{1}{2(\alpha+1)}} \quad (39)$$

if  $n^{\frac{1}{2(\alpha+1)}} \leq p$ , which is different from the choice of  $k$  for the operator norm in equation (15). If  $n^{\frac{1}{2(\alpha+1)}} > p$ , we will choose  $k = p$ , then the bias part is 0 and consequently

$$\mathbb{E} \frac{1}{p} \left\| \check{\Sigma} - \Sigma \right\|_F^2 \leq C \frac{p}{n}. \quad \blacksquare$$

**Remark 4** Under the Frobenius norm the optimal tapering parameter  $k$  is of the order  $n^{\frac{1}{2(\alpha+1)}}$ . The rate of convergence of the tapering estimator with  $k \asymp n^{\frac{1}{2(\alpha+1)}}$  under the operator norm is

$$\frac{\log p}{n} + n^{-\frac{\alpha}{\alpha+1}},$$

which is slower than  $n^{-\frac{2\alpha}{2\alpha+1}} + \frac{\log p}{n}$  in equation (1). Similarly the optimal procedure under the operator norm is not rate optimal under the Frobenius norm. Therefore, the optimal choice of the tapering parameter  $k$  critically depends on the norm under which the estimation accuracy is measured.

**Remark 5** Similarly it can be shown that under the Frobenius norm the banding estimator with  $k \asymp n^{\frac{1}{2(\alpha+1)}}$  is rate optimal. Under the operator norm, Bickel and Levina (2008a) chose  $k \asymp \left(\frac{\log p}{n}\right)^{\frac{1}{2(\alpha+1)}}$  for the banding estimator which is close to  $n^{\frac{1}{2(\alpha+1)}}$  up to a logarithmic factor of  $p$ .

**Remark 6** For the parameter space  $\mathcal{F}_\alpha$  given in (3), we have

$$\frac{1}{p} \sum_{\{(i,j):k<|i-j|\}} \sigma_{ij}^2 \leq \frac{1}{p} \left( \sum_{\{(i,j):k<|i-j|\}} |\sigma_{ij}| \right)^2 \leq Ck^{-2\alpha},$$

which is different from the upper bound in equation (37). As a consequence, under the Frobenius norm, the optimal rates for estimating  $\Sigma$  over  $\mathcal{P}_\alpha$  and over  $\mathcal{P}'_\alpha$  are different. In contrast, as seen in Section 3, the minimax rates over the two parameter spaces under the operator norm are the same.

## 4.2 Lower Bound under the Frobenius Norm

We only establish the lower bound for the parameter space  $\mathcal{P}'_\alpha$  given in (30). Again the argument for  $\mathcal{P}_\alpha$  is similar. As in the case of estimation under the operator norm, we need to construct a finite collection of multivariate normal distributions with a parameter space  $\mathcal{G}_2 \subset \mathcal{G}_\alpha$  such that

$$\inf_{\hat{\Sigma}} \sup_{\mathcal{G}_2} \mathbb{E} \frac{1}{p} \left\| \hat{\Sigma} - \Sigma \right\|_F^2 \geq c \frac{k}{n}.$$

for some  $c > 0$  when  $k = \min \left\{ n^{\frac{1}{2(\alpha+1)}}, p/2 \right\}$ .

We construct  $\mathcal{G}_2$  as follows. Let  $0 < \tau < M$  be a constant. Define

$$\mathcal{G}_2 = \left\{ \Sigma(\theta) : \Sigma(\theta) = I + \left( \theta_{ij} \tau n^{-\frac{1}{2}} I \{1 \leq |i-j| \leq k\} \right)_{p \times p}, \text{ for } \theta_{ij} = \theta_{ji} = 0 \text{ or } 1 \right\}$$

It is easy to verify that  $\mathcal{G}_2 \subset \mathcal{G}_\alpha$  as  $n \rightarrow \infty$ . Note that  $\theta \in \Theta = \{0, 1\}^{kp-k(k+1)/2}$ .

Applying Assouad's Lemma with  $d$  the Frobenius norm and  $s = 2$  to the parameter space  $\mathcal{G}_2$ , we have

$$\max_{\theta \in \mathcal{G}_2} 2^2 E_\theta \frac{1}{p} \left\| \hat{\Sigma} - \Sigma(\theta) \right\|_F^2 \geq \min_{H(\theta, \theta') \geq 1} \frac{\frac{1}{p} \left\| \Sigma(\theta) - \Sigma(\theta') \right\|_F^2}{H(\theta, \theta')} \frac{kp - k(k+1)/2}{2} \min_{H(\theta, \theta')=1} \|P_\theta \wedge P_{\theta'}\|.$$

Note that

$$\min_{H(\theta, \theta') \geq 1} \frac{1}{p} \frac{\|\Sigma(\theta) - \Sigma(\theta')\|_F^2}{H(\theta, \theta')} = \min_{H(\theta, \theta') \geq 1} \frac{1}{p} \frac{\left[\tau n^{-\frac{1}{2}}\right]^2 \sum |\theta_{ij} - \theta'_{ij}|^2}{H(\theta, \theta')} = \frac{\tau^2}{p} n^{-1}.$$

It is easy to see that

$$\frac{kp - k(k+1)/2}{2} \asymp kp.$$

**Lemma 8** *Let  $P_\theta$  be the joint distribution of  $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{iid}{\sim} N(0, \Sigma(\theta))$  with  $\Sigma(\theta) \in \mathcal{G}_2$ . Then for some constant  $c_1 > 0$  we have*

$$\min_{H(\theta, \theta')=1} \|P_\theta \wedge P_{\theta'}\| \geq c_1.$$

We omit the proof of this lemma. It is very similar to and simpler than the proof of Lemma 6.

From Lemma 8 we have for some  $c > 0$

$$\min_{H(\theta, \theta')=1} \|P_\theta \wedge P_{\theta'}\| \geq c \tag{40}$$

thus

$$\max_{\theta \in \mathcal{G}_2} 2^2 E_\theta \frac{1}{p} \left\| \hat{\Sigma} - \Sigma(\theta) \right\|_F^2 \geq c \min \left\{ n^{-\frac{2\alpha+1}{2(\alpha+1)}}, \frac{p}{n} \right\}$$

which implies that the rate obtained in (38) is optimal. ■

## 5 Estimation of the Inverse Covariance Matrix

The inverse of the covariance matrix  $\Sigma^{-1}$  is of significant interest in many statistical applications. The results and analysis given in Section 3 can be used to derive the optimal rate of convergence for estimating  $\Sigma^{-1}$  under the operator norm.

For estimating the inverse covariance matrix  $\Sigma^{-1}$  we require the minimum eigenvalue of  $\Sigma$  to be bounded away from zero. For  $\delta > 0$ , we define

$$L_\delta = \{\Sigma : \lambda_{\min}(\Sigma) \geq \delta\}. \tag{41}$$

Let  $\tilde{\mathcal{P}}_\alpha = \tilde{\mathcal{P}}_\alpha(M_0, M, \rho, \delta)$  denote the set of distributions of  $\mathbf{X}_1$  that satisfy (3), (7) and (41), and similarly distributions in  $\tilde{\mathcal{P}}'_\alpha = \tilde{\mathcal{P}}'_\alpha(M_0, M, \rho, \delta)$  satisfy (7), (30) and (41).

The following theorem gives the minimax rate of convergence for estimating  $\Sigma^{-1}$ .

**Theorem 5** *The minimax risk of estimating the inverse covariance matrix  $\Sigma^{-1}$  satisfies*

$$\inf_{\hat{\Sigma}} \sup_{\tilde{\mathcal{P}}} \mathbb{E} \left\| \hat{\Sigma}^{-1} - \Sigma^{-1} \right\|^2 \asymp \min \left\{ n^{-\frac{2\alpha}{2\alpha+1}} + \frac{\log p}{n}, \frac{p}{n} \right\} \quad (42)$$

where  $\tilde{\mathcal{P}}$  denotes either  $\tilde{\mathcal{P}}_\alpha$  or  $\tilde{\mathcal{P}}'_\alpha$ .

*Proof of Theorem 5:* We shall focus on the case  $p \geq n^{\frac{1}{2\alpha+1}}$ . The proof for the case of  $p < n^{\frac{1}{2\alpha+1}}$  is similar. To establish the upper bound, note that

$$\hat{\Sigma}^{-1} - \Sigma^{-1} = \hat{\Sigma}^{-1} (\Sigma - \hat{\Sigma}) \Sigma^{-1}$$

then

$$\left\| \hat{\Sigma}^{-1} - \Sigma^{-1} \right\|^2 = \left\| \hat{\Sigma}^{-1} (\Sigma - \hat{\Sigma}) \Sigma^{-1} \right\|^2 \leq \left\| \hat{\Sigma}^{-1} \right\|^2 \left\| \Sigma - \hat{\Sigma} \right\|^2 \left\| \Sigma^{-1} \right\|^2.$$

It follows from Assumption (3) that  $\left\| \Sigma^{-1} \right\|^2 \leq \delta^{-2}$ . Note that  $\mathbb{P} \left\{ \left\| \check{\Sigma} - \mathbb{E} \check{\Sigma} \right\|^2 > \epsilon \right\} \leq 4p5^m \exp(-n\epsilon^2\rho_1)$  for any  $\epsilon > 0$  which decays faster than any polynomial of  $n$  as shown in the proof of Lemmas 2 and 3. Let  $\lambda_{\min}(\check{\Sigma})$  and  $\lambda_{\min}(\mathbb{E}\check{\Sigma})$  be the smallest eigenvalues of  $\check{\Sigma}$  and  $\mathbb{E}\check{\Sigma}$  respectively. Then  $\mathbb{P} \left( \lambda_{\min}(\check{\Sigma}) \leq \lambda_{\min}(\mathbb{E}\check{\Sigma}) - \epsilon^{1/2} \right) \geq \mathbb{P} \left( \left| \lambda_{\min}(\check{\Sigma}) - \lambda_{\min}(\mathbb{E}\check{\Sigma}) \right| \geq \epsilon^{1/2} \right)$  decays faster than any polynomial of  $n$ . Let  $0 < \epsilon < \left[ \lambda_{\min}(\mathbb{E}\check{\Sigma}) / 2 \right]^2$  and  $c = 1 / \left[ \lambda_{\min}(\mathbb{E}\check{\Sigma}) - \epsilon^{1/2} \right]$ , then  $\mathbb{P} \left( \left\| \hat{\Sigma}^{-1} \right\| \geq c \right)$  decays faster than any polynomial of  $n$ . Therefore,

$$\begin{aligned} \mathbb{E} \left\| \hat{\Sigma}^{-1} - \Sigma^{-1} \right\|^2 &\leq \left( \frac{c}{\delta} \right)^2 \mathbb{E} \left\| \Sigma - \hat{\Sigma} \right\|^2 + \mathbb{E} \left[ \left\| \hat{\Sigma}^{-1} \right\|^2 \left\| \Sigma - \hat{\Sigma} \right\|^2 \left\| \Sigma^{-1} \right\|^2 I \left( \left\| \hat{\Sigma}^{-1} \right\| \geq c \right) \right] \\ &\leq C \min \left\{ n^{-\frac{2\alpha}{2\alpha+1}} + \frac{\log p}{n}, \frac{p}{n} \right\}. \end{aligned}$$

The proof of the lower bound is almost identical to that of Theorem 1 except that here we need to show

$$\min_{H(\theta, \theta') \geq 1} \frac{\left\| \Sigma^{-1}(\theta) - \Sigma^{-1}(\theta') \right\|^2}{H(\theta, \theta')} \geq cka^2$$

instead of Lemma 5. For a positive definite matrix  $A$ , let  $\lambda_{\min}(A)$  denote the minimum eigenvalue of  $A$ . Since

$$\Sigma^{-1}(\theta) - \Sigma^{-1}(\theta') = \Sigma^{-1}(\theta') (\Sigma(\theta) - \Sigma(\theta')) \Sigma^{-1}(\theta),$$

we have

$$\left\| \Sigma^{-1}(\theta) - \Sigma^{-1}(\theta') \right\| \geq \lambda_{\min}(\Sigma^{-1}(\theta)) \lambda_{\min}(\Sigma^{-1}(\theta')) \left\| \Sigma(\theta) - \Sigma(\theta') \right\|.$$

Note that

$$\lambda_{\min}(\Sigma^{-1}(\theta)) > 1/M_0, \lambda_{\min}(\Sigma^{-1}(\theta')) > 1/M_0,$$

then Lemma 5 implies

$$\min_{H(\theta, \theta') \geq 1} \frac{\|\Sigma^{-1}(\theta) - \Sigma^{-1}(\theta')\|^2}{H(\theta, \theta')} \geq M_0^{-4} \min_{H(\theta, \theta') \geq 1} \frac{\|\Sigma(\theta) - \Sigma(\theta')\|^2}{H(\theta, \theta')} \geq cka^2$$

for some constant  $c > 0$ . ■

## 6 Simulation Study

We now turn to the numerical performance of the proposed tapering estimator and compare it with that of the banding estimator of Bickel and Levina (2008a). In the numerical study, we shall consider estimating a covariance matrix in the parameter space  $\mathcal{F}_\alpha$  defined in (3). Specifically, we consider the covariance matrix  $\Sigma = (\sigma_{ij})_{1 \leq i, j \leq p}$  of the form

$$\sigma_{ij} = \begin{cases} 1, & 1 \leq i = j \leq p \\ \rho|i - j|^{-(\alpha+1)}, & 1 \leq i \neq j \leq p \end{cases}. \quad (43)$$

Note that this matrix is a Toeplitz matrix. But we do not assume that the structure is known and do not use the information in any estimation procedure.

The banding estimator in (29) depends on the choice of  $k$ . An optimal tradeoff of  $k$  is  $k \asymp (n/\log p)^{1/(2\alpha+2)}$  as discussed in Section 3.3. See Bickel and Levina (2006). The tapering estimator (6) also depends on  $k$  for which the optimal tradeoff is  $k \asymp n^{1/(2\alpha+1)}$ . In our simulation study, we choose  $k = \lfloor (n/\log p)^{1/(2\alpha+2)} \rfloor$  for the banding estimator, and  $k = \lfloor n^{1/(2\alpha+1)} \rfloor$  for the tapering estimator.

A range of parameter values for  $\alpha$ ,  $n$  and  $p$  are considered. Specifically,  $\alpha$  ranges from 0.1 to 0.5, the sample size  $n$  ranges from 250 to 3000 and the dimension  $p$  goes from 250 to 3000. We choose the value of  $\rho$  to be  $\rho = 0.6$  so that all matrices are non-negative definite and their smallest eigenvalues are close to 0. Table 1 reports the average errors under the spectral norm over 100 replications for the two procedures. The cases where the tapering estimator underperforms the banding estimator are highlighted in boldface. Figure 2 plots the ratios of the average errors of the banding estimator to the corresponding average errors of the tapering estimator for  $\alpha = 0.1, 0.2, 0.3$  and 0.5. The case of  $\alpha = 0.4$  is similar to the case of  $\alpha = 0.3$ .

It can be seen from Table 1 and figure 2 that the tapering estimator outperforms the banding estimator in 121 out of 125 cases. For the given dimension  $p$ , the ratio of the

average error of the banding estimator to the corresponding average error of the tapering estimator tends to increase as the sample size  $n$  increases. The tapering estimator fails to outperform the banding estimator only when  $\alpha = 0.5$  and  $n = 250$  in which case the values of  $k$  are small for both estimators.

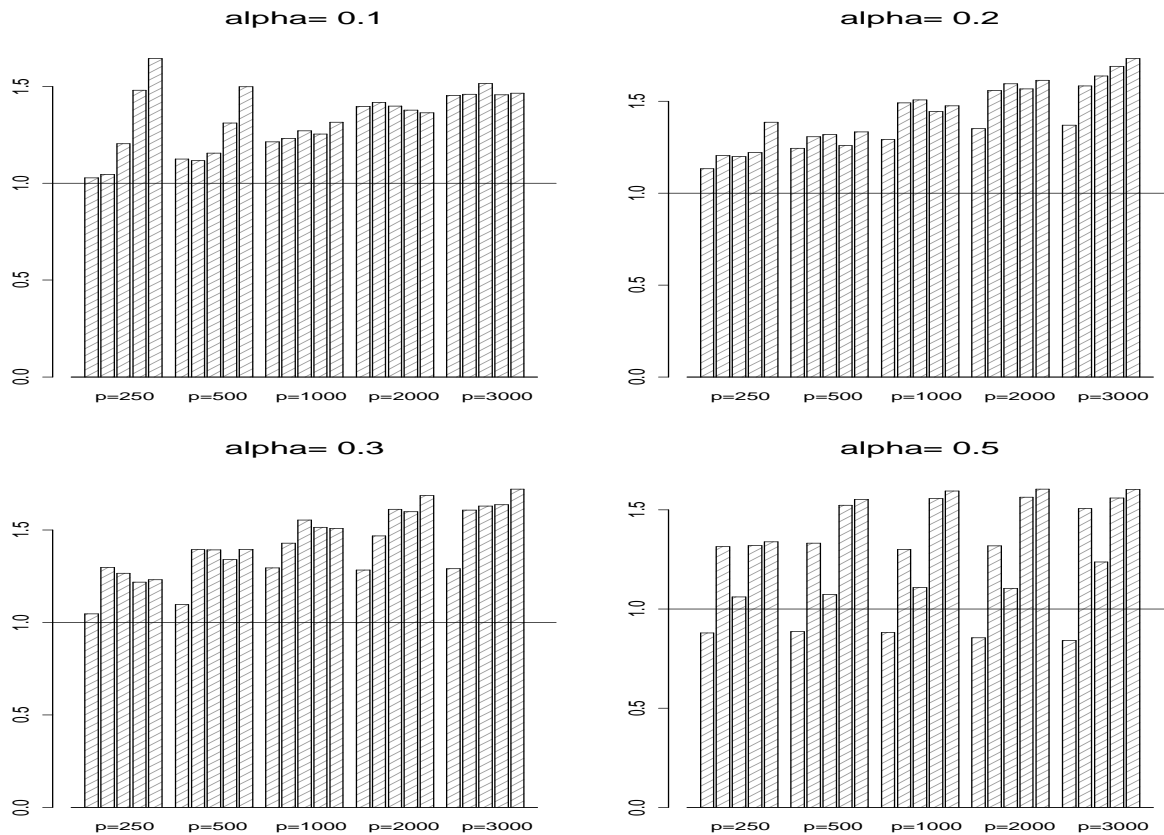


Figure 2: The vertical bars represent the ratios of the average error of the banding estimator to the corresponding average error of the tapering estimator. The higher the bar the better the relative performance of the tapering estimator. For each value of  $p$  the bars are ordered from left to right by the sample sizes ( $n = 250$  to 3000).

**Remark 7** We have also carried out additional simulations for larger values of  $\alpha$  with the same sample sizes and dimensions. The performance of the tapering and banding estimators are similar. This is mainly due to the fact that the values of  $k$  for both estimators are very small for large  $\alpha$  when  $n$  and  $p$  are only moderately large.

$p$	$n$	$\alpha = 0.1$		$\alpha = 0.2$		$\alpha = 0.3$		$\alpha = 0.4$		$\alpha = 0.5$	
		BL	CZZ	BL	CZZ	BL	CZZ	BL	CZZ	BL	CZZ
250	250	2.781	2.706	2.291	2.023	1.762	1.684	1.618	1.517	<b>1.325</b>	<b>1.507</b>
	500	2.409	2.302	1.898	1.575	1.562	1.204	1.361	1.185	1.080	0.822
	1000	2.029	1.685	1.631	1.361	1.289	1.018	1.056	0.795	0.911	0.859
	2000	1.706	1.153	1.369	1.122	1.106	0.908	0.878	0.655	0.715	0.542
	3000	1.522	0.926	1.242	0.896	0.983	0.798	0.810	0.658	0.645	0.482
500	250	3.277	2.914	2.609	2.097	1.961	1.788	1.745	1.610	<b>1.392</b>	<b>1.571</b>
	500	2.901	2.598	2.199	1.683	1.751	1.256	1.475	1.234	1.152	0.865
	1000	2.539	2.197	1.942	1.472	1.481	1.064	1.178	0.843	0.984	0.917
	2000	2.263	1.726	1.669	1.326	1.293	0.965	1.067	0.700	0.866	0.569
	3000	2.066	1.379	1.538	1.154	1.220	0.874	0.919	0.696	0.781	0.503
1000	250	3.747	3.086	2.873	2.223	2.385	1.842	1.833	1.694	<b>1.449</b>	<b>1.643</b>
	500	3.370	2.735	2.635	1.768	1.906	1.334	1.565	1.297	1.203	0.925
	1000	3.097	2.437	2.315	1.536	1.741	1.121	1.382	0.883	1.037	0.936
	2000	2.730	2.177	2.011	1.392	1.523	1.006	1.156	0.722	0.920	0.591
	3000	2.589	1.968	1.865	1.264	1.374	0.911	1.072	0.723	0.834	0.523
2000	250	4.438	3.177	3.107	2.300	2.511	1.956	1.903	1.744	<b>1.484</b>	<b>1.736</b>
	500	3.969	2.800	2.868	1.841	2.030	1.383	1.638	1.356	1.239	0.940
	1000	3.538	2.531	2.551	1.599	1.866	1.158	1.452	0.912	1.074	0.973
	2000	3.242	2.353	2.248	1.434	1.649	1.031	1.224	0.751	0.955	0.611
	3000	3.025	2.219	2.101	1.302	1.566	0.929	1.141	0.743	0.868	0.541
3000	250	4.679	3.219	3.230	2.358	2.576	1.995	1.931	1.797	<b>1.494</b>	<b>1.776</b>
	500	4.214	2.887	2.991	1.890	2.282	1.419	1.664	1.384	1.463	0.971
	1000	3.901	2.575	2.674	1.633	1.933	1.186	1.482	0.929	1.224	0.990
	2000	3.488	2.395	2.452	1.451	1.717	1.049	1.254	0.768	0.965	0.619
	3000	3.336	2.278	2.288	1.321	1.632	0.948	1.172	0.750	0.880	0.549

Table 1: The average errors under the spectral norm of the banding estimator (BL) and the tapering estimator (CZZ) over 100 replications. The cases where the tapering estimator underperforms the banding estimator are highlighted in boldface.

## 7 Proofs of Auxiliary Lemmas

In this section we give proofs of auxiliary lemmas stated and used in Sections 3 - 5.

**Proof of Lemma 1:** Without loss of generality we assume that  $i \leq j$ . The set  $\{i, j\}$  is contained in the set  $\{l, \dots, l + k_h - 1\}$  if and only if  $l \leq i \leq j \leq l + k_h - 1$ , i.e.,  $j - k_h + 1 \leq l \leq i$ . Note that  $\text{Card}\{l : j - k_h + 1 \leq l \leq i\} = (i - (j - k_h + 1) + 1)_+ = (k_h - |i - j|)_+$ , then  $\text{Card}\{l : \{i, j\} \subset \{l, \dots, l + k_h - 1\}\} = (k_h - |i - j|)_+$ . Similarly we have  $\text{Card}\{l : \{i, j\} \subset \{l, \dots, l + k - 1\}\} = (k - |i - j|)_+$ . Thus we have

$$\begin{aligned} kw_{ij} &= (k - |i - j|)_+ - (k_h - |i - j|)_+ \\ &= \text{Card}\{l : \{i, j\} \subset \{l, \dots, l + k - 1\}\} - \text{Card}\{l : \{i, j\} \subset \{l, \dots, l + k_h - 1\}\}. \blacksquare \end{aligned}$$

**Proof of Lemma 2:** Without loss of generality we assume that  $p$  is divisible by  $m$ . Recall that  $M_l^{(m)} = (\tilde{\sigma}_{ij} I \{l \leq i < l + m, l \leq j < l + m\})_{p \times p}$ . Note that  $M_l^{(m)}$  is empty when  $l \leq 1 - m$ , and has at least one nonzero entry when  $l \geq 2 - m$ . Set  $\delta_l^{(m)} = M_l^{(m)} - \mathbb{E}M_l^{(m)}$  and  $S^{(m)} = \sum_{l=2-m}^p M_l^{(m)}$ . It follows from (6) that

$$\left\| S^{(m)} - \mathbb{E}S^{(m)} \right\| \leq \sum_{l=1}^m \left\| \sum_{-1 \leq j < p/m} \delta_{jm+l}^{(m)} \right\|. \quad (44)$$

Since  $\delta_{jm+l}^{(m)}$  are disjoint diagonal blocks over  $-1 \leq j < p/m$ , we have

$$\left\| S^{(m)} - \mathbb{E}S^{(m)} \right\| \leq m \max_{1 \leq l \leq m} \left\| \sum_{-1 \leq j < p/m} \delta_{jm+l}^{(m)} \right\| \leq m \max_{1-m \leq l \leq p} \left\| \delta_l^{(m)} \right\|. \quad (45)$$

Since  $\delta_l^{(k_h)}$  and  $\delta_l^{(k)}$  are all sub-blocks of certain matrix  $\delta_l^{(k)}$  with  $1 \leq l \leq p - k + 1$ , Lemma 2 now follows immediately from equations (45) and (6).  $\blacksquare$

**Proof of Lemma 3:** For any  $m \times m$  symmetric matrix  $A$ , we have

$$|u^T A u| - |v^T A v| \leq |u^T A u - v^T A v| = \left| (u - v)^T A (u + v) \right| \leq \|u - v\| \|A\| \|u + v\|$$

Let  $S_{1/2}^{m-1}$  be a  $1/2$  net of the unit sphere  $S^{m-1}$  in the Euclidean distance in  $\mathbb{R}^m$ . We have

$$\|A\| \leq \sup_{u \in S^{m-1}} |u^T A u| \leq \sup_{u \in S_{1/2}^{m-1}} |u^T A u| + \frac{1}{2} \|A\| \frac{3}{2} = \sup_{u \in S_{1/2}^{m-1}} |u^T A u| + \frac{3}{4} \|A\|$$

which implies  $\|A\| \leq 4 \sup_{u \in S_{1/2}^{m-1}} |u^T A u|$ . Since we are allowed to pack  $\text{Card}(S_{1/2}^{m-1})$  balls of radius  $1/4$  into a  $1 + 1/4$  ball in  $\mathbb{R}^m$ , volume comparison yields

$$(1/4)^m \text{Card}(S_{1/2}^{m-1}) \leq (5/4)^m,$$

i.e.,  $\text{Card}(S_{1/2}^{m-1}) \leq 5^m$ . Thus there exist  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{5^m} \in S^{m-1}$  such that

$$\|A\| \leq 4 \sup_{j \leq 5^m} |v_j^T A v_j|, \text{ for all } m \times m \text{ symmetric } A.$$

This one step approximation argument is similar to the proof of Proposition 4.2 (ii) in Zhang and Huang (2008).

Let  $\mathbf{X}_1, \dots, \mathbf{X}_n$  be i.i.d.  $p$ -vectors with  $\mathbb{E}(\mathbf{X}_1 - \mu)(\mathbf{X}_1 - \mu)^T = \Sigma$ . Under the subgaussian assumption in (7) there exists  $\rho > 0$  such that

$$\mathbb{P}\{\mathbf{v}^T(\mathbf{X}_i - \mathbb{E}\mathbf{X}_i)(\mathbf{X}_i - \mathbb{E}\mathbf{X}_i)^T \mathbf{v} > x\} \leq e^{-x\rho/2} \text{ for all } x > 0 \text{ and } \|\mathbf{v}\| = 1$$

which implies  $\mathbb{E}(t\mathbf{v}^T(\mathbf{X}_i - \mathbb{E}\mathbf{X}_i)(\mathbf{X}_i - \mathbb{E}\mathbf{X}_i)^T \mathbf{v}) < \infty$  for all  $t < \rho/2$  and  $\|\mathbf{v}\| = 1$ , then there exists  $\rho_1 > 0$  such that

$$\mathbb{P}\left\{\left|\frac{1}{n} \sum_{i=1}^n \mathbf{v}^T [(\mathbf{X}_i - \mathbb{E}\mathbf{X}_i)(\mathbf{X}_i - \mathbb{E}\mathbf{X}_i)^T - \Sigma] \mathbf{v}\right| > x\right\} \leq e^{-nx^2\rho_1/2}$$

for all  $0 < x < \rho_1$  and  $\|\mathbf{v}\| = 1$ . (See, e.g., Chapter 2 in Saulis and Statulevicius (1991).)

Thus we have

$$\begin{aligned} \mathbb{P}\left\{\max_{1 \leq l \leq p-m+1} \|M_l^{(m)} - \mathbb{E}M_l^{(m)}\| > x\right\} &\leq \sum_{1 \leq l \leq p-m+1} \mathbb{P}\left\{\|M_l^{(m)} - \mathbb{E}M_l^{(m)}\| > x\right\} \\ &\leq 2p5^m \sup_{\mathbf{v}_j, l} \mathbb{P}\{|\mathbf{v}_j^T (M_l^{(m)} - \mathbb{E}M_l^{(m)}) \mathbf{v}_j| > x\} \\ &\leq 2p5^m \exp(-nx^2\rho_1/2). \quad \blacksquare \end{aligned}$$

**Proof of Lemma 5:** Set  $v = (1 \{k_h \leq i \leq k\})$  and let

$$(w_i) = [\Sigma(\theta) - \Sigma(\theta')] v.$$

Note that there are exactly  $H(\theta, \theta')$  number of  $w_i$  such that  $|w_i| = \tau k_h a$ , and  $\|v\|_2^2 = k_h$ . This implies

$$\|\Sigma(\theta) - \Sigma(\theta')\|^2 \geq \frac{\|[\Sigma(\theta) - \Sigma(\theta')] v\|_2^2}{\|v\|_2^2} \geq \frac{H(\theta, \theta') \cdot (\tau k a)^2}{k_h} = H(\theta, \theta') \cdot \tau^2 k_h a^2. \quad \blacksquare$$

**Proof of Lemma 6:** When  $H(\theta, \theta') = 1$ , we will show

$$\begin{aligned} \|P_{\theta'} - P_{\theta}\|_1^2 &\leq 2K(P_{\theta'}|P_{\theta}) = 2n \left[ \frac{1}{2} \text{tr}(\Sigma(\theta') \Sigma^{-1}(\theta)) - \frac{1}{2} \log \det(\Sigma(\theta') \Sigma^{-1}(\theta)) - \frac{p}{2} \right] \\ &\leq n \cdot cka^2 \end{aligned}$$

for some small  $c > 0$ , where  $K(\cdot|\cdot)$  is the Kullback–Leibler divergence and the first inequality follows from the well known Pinsker’s inequality (see, e.g., Csiszár (1967)). This immediately implies the  $L_1$  distance between two measures is bounded away from 1, and then the lemma follows. Write

$$\Sigma(\theta') = D_1 + \Sigma(\theta).$$

Then

$$\frac{1}{2} \text{tr}(\Sigma(\theta') \Sigma^{-1}(\theta)) - \frac{p}{2} = \frac{1}{2} \text{tr}(D_1 \Sigma^{-1}(\theta)).$$

Let  $\lambda_i$  be the eigenvalues of  $D_1 \Sigma^{-1}(\theta)$ . Since  $D_1 \Sigma^{-1}(\theta)$  is similar to the symmetric matrix  $\Sigma^{-1/2}(\theta) D_1 \Sigma^{-1/2}(\theta)$ , and

$$\left\| \Sigma^{-1/2}(\theta) D_1 \Sigma^{-1/2}(\theta) \right\| \leq \left\| \Sigma^{-1/2}(\theta) \right\| \|D_1\| \left\| \Sigma^{-1/2}(\theta) \right\| \leq c_1 \|D_1\| \leq c_1 \|D_1\|_1 \leq c_2 ka,$$

then all eigenvalues  $\lambda_i$ ’s are real and in the interval  $[-c_2 ka, c_2 ka]$ , where  $ka = k \cdot k^{-(\alpha+1)} = k^{-\alpha} \rightarrow 0$ . Note that the Taylor expansion yields

$$\log \det(\Sigma(\theta') \Sigma^{-1}(\theta)) = \log \det(I + D_1 \Sigma^{-1}(\theta)) = \text{tr}(D_1 \Sigma^{-1}(\theta)) - R_3$$

where

$$R_3 \leq c_3 \sum_{i=1}^p \lambda_i^2 \text{ for some } c_3 > 0.$$

Write  $\Sigma^{-1/2}(\theta) = UV^{1/2}U^T$ , where  $UU^T = I$  and  $V$  is a diagonal matrix. It follows from the fact that the Frobenius norm of a matrix remains the same after an orthogonal transformation that

$$\sum_{i=1}^p \lambda_i^2 = \left\| \Sigma^{-1/2}(\theta) D_1 \Sigma^{-1/2}(\theta) \right\|_F^2 \leq \|V\|^2 \cdot \|U^T D_1 U\|_F^2 = \|\Sigma^{-1}(\theta)\|^2 \cdot \|D_1\|_F^2 \leq c_4 ka^2. \quad \blacksquare$$

**Acknowledgment:** The authors would like to thank James X. Hu for assistance in carrying out the simulation study in Section 6. We also thank the Associate Editor and three referees for thorough and useful comments which have helped to improve the presentation of the paper.

## References

- [1] Assouad, P. (1983). Deux remarques sur l'estimation, *C. R. Acad. Sci. Paris* **296**, 1021-1024
- [2] Bickel, P. J. and Levina, E. (2008a). Regularized estimation of large covariance matrices. *Ann. Statist.* **36**, 199-227.
- [3] Bickel, P. J. and Levina, E. (2008b). Covariance regularization by thresholding. *Ann. Statist.* **36**, 2577-2604.
- [4] Csiszár, I. (1967). Information-type measures of difference of probability distributions and indirect observation, *Studia Sci. Math. Hungar.* **2**, 229-318.
- [5] Davidson, K. R. and Szarek, S. J. (2001). Local operator theory, random matrices and Banach spaces, *Handbook in Banach Spaces* Vol I, ed. W. B. Johnson, J. Lindenstrauss, Elsevier, 317-366.
- [6] El Karoui, N. (2008). Operator norm consistent estimation of large dimensional sparse covariance matrices. *Ann. Statist.* **36**, 2717-2756.
- [7] Golub, G. H. and Van Loan, C. F. (1983). *Matrix Computations*. The John Hopkins University Press, Baltimore, Maryland.
- [8] Fan, J., Fan, Y., and Lv, J. (2006). High dimensional covariance matrix estimation using a factor model. *Journal of Econometrics*, to appear.
- [9] Furrer, R. and Bengtsson, T. (2007). Estimation of high-dimensional prior and posterior covariance matrices in Kalman filter variants, *Journal of Multivariate Analysis* **98**, 227-255.
- [10] Huang, J., Liu, N., Pourahmadi, M. and Liu, L. (2006), Covariance matrix selection and estimation via penalised normal likelihood, *Biometrika* **93**, 85-98.
- [11] Johnstone, I. M. (2001). On the distribution of the largest eigenvalue in principal components analysis, *Ann. Statist.* **29**, 295-327.
- [12] Johnstone, I. M. and Lu, A. Y. (2004). Sparse principal components analysis, *J. Amer. Statist. Assoc.*, tentatively accepted.
- [13] Lam, C. and Fan, J. (2007). Sparsistency and rates of convergence in large covariance matrices estimation. Princeton University. Technical report.

- [14] Meinshausen, N. and Bühlmann, P. (2006). High dimensional graphs and variable selection with the Lasso. *Ann. Statist.* **34**, 1436-1462.
- [15] Muirhead, R. J. (1987). Developments in eigenvalue estimation. *Advances in Multivariate Statistical*, (A. K. Gupta, Ed.), 277-288, Reidel, Boston.
- [16] Ravikumar, P., Wainwright, M. J., Raskutti, G., and Yu, B. (2008). High-dimensional covariance estimation by minimizing  $l_1$ -penalized log-determinant divergence. Technical Report, UC Berkeley.
- [17] Rothman, A.J. , Bickel, P. J., Levina, E. and Zhu, J. (2008). Sparse permutation invariant covariance estimation. *Electronic Journal of Statistics* **2**, 494-515.
- [18] Rudelson, M. and Vershynin, R. (2007). Sampling from large matrices: an approach through geometric functional analysis, *Journal of the ACM* **54**, No. 4, Article 21.
- [19] Saulis, L and Statulevicius, V.A. (1991). *Limit Theorems For Large Deviations*. Kluwer Academic Publishers.
- [20] Wu, W. B. and Pourahmadi, M. (2009). Banding Sample Covariance Matrices of Stationary Processes. *Statistica Sinica*. To appear.
- [21] Yu, B. (1997). Assouad, Fano, and Le Cam. *Festschrift for Lucien Le Cam*. D. Pollard, E. Torgersen, and G. Yang (eds), pp. 423-435, Springer-Verlag.
- [22] Zhang, C.-H. and Huang, J. (2008). The sparsity and bias of the Lasso selection in high-dimensional linear regression. *Ann. Statist.* **36**, 1567–1594.
- [23] Zou, H., Hastie, T., and Tibshirani, R. (2006). Sparse principal components analysis. *Journal of Computational and Graphical Statistics* **15**, 265-286.