



Published in final edited form as:

*J Am Stat Assoc.* 2018 ; 113(521): 380–389. doi:10.1080/01621459.2016.1256815.

## Embracing the Blessing of Dimensionality in Factor Models

Quefeng Li<sup>a</sup>, Guang Cheng<sup>b</sup>, Jianqing Fan<sup>c,d</sup>, and Yuyan Wang<sup>c</sup>

<sup>a</sup>Department of Biostatistics, University of North Carolina at Chapel Hill, NC

<sup>b</sup>Department of Statistics, Purdue University, IN

<sup>c</sup>Department of Operations Research and Financial Engineering, Princeton University, NJ

<sup>d</sup>School of Data Science, Fudan University, Shanghai, China

### Abstract

Factor modeling is an essential tool for exploring intrinsic dependence structures among high-dimensional random variables. Much progress has been made for estimating the covariance matrix from a high-dimensional factor model. However, the blessing of dimensionality has not yet been fully embraced in the literature: much of the available data are often ignored in constructing covariance matrix estimates. If our goal is to accurately estimate a covariance matrix of a set of targeted variables, shall we employ additional data, which are beyond the variables of interest, in the estimation? In this article, we provide sufficient conditions for an affirmative answer, and further quantify its gain in terms of Fisher information and convergence rate. In fact, even an oracle-like result (as if all the factors were known) can be achieved when a sufficiently large number of variables is used. The idea of using data as much as possible brings computational challenges. A divide-and-conquer algorithm is thus proposed to alleviate the computational burden, and also shown not to sacrifice any statistical accuracy in comparison with a pooled analysis. Simulation studies further confirm our advocacy for the use of full data, and demonstrate the effectiveness of the above algorithm. Our proposal is applied to a microarray data example that shows empirical benefits of using more data. Supplementary materials for this article are available online.

### Keywords

Asymptotic normality; Auxiliary data; Divide-and-conquer; Factor model; Fisher information; High-dimensionality

## 1. Introduction

With the advance of modern information technology, it is now possible to track millions of variables or subjects simultaneously. To discover the relationship among them, the estimation of a high-dimensional covariance matrix  $\Sigma$  has recently received a great deal of attention in the literature. Researchers proposed various regularization methods to obtain

consistent estimators of  $\Sigma$  (Bickel and Levina 2008; Rothman et al. 2008; Lam and Fan 2009; Cai, Zhang, and Zhou 2010; Cai and Liu 2011). A key assumption for these regularization methods is that  $\Sigma$  is sparse, that is, many elements of  $\Sigma$  are small or exactly zero.

Different from such a sparsity condition, factor analysis assumes that the intrinsic dependence is mainly driven by some common latent factors (Johnson and Wichern 1992). For example, in modeling stock returns, Fama and French (1993) proposed the well-known Fama–French three-factor model. In the factor model,  $\Sigma$  has spiked eigenvalues and dense entries. In the high-dimensional setting, there are many recent studies on the estimation of the covariance matrix based on the factor model (Fan, Fan, and Lv 2008; Fan, Liao, and Mincheva 2011, 2013; Bai and Li 2012; Bai and Liao 2013), where the number of variables can be much larger than the number of observations.

The interest of this article is on the estimation of the covariance matrix for a certain set of variables using auxiliary data information. In the literature, we use only the data information on the variables of interest. In the data-rich environment today, substantially more amount of data information is indeed available, but is often ignored in statistical analysis. For example, we might be interested in understanding the covariance matrix of 50 stocks in a portfolio, yet the available data information is a time series of thousands of stocks. Similarly, an oncologist may wish to study the dependence or network structures among 100 genes that are significantly associated with a certain cancer, yet she has expression data for over 20,000 genes from the whole genome. Can we benefit from using much more rich auxiliary data?

The answer to the above question is affirmative when a factor model is imposed. Since the whole system is driven by a few common factors, these common factors can be inferred more accurately from a much larger set of data information (Fan, Liao, and Mincheva 2013), which is indeed a “blessing of dimensionality.” A major contribution of this article is to characterize how much the estimation of the covariance matrix of interest and also common factors can be improved by auxiliary data information (and under what conditions).

Consider the following factor model for all  $p$  observable data  $\mathbf{y}_t = (y_{1t}, \dots, y_{pt})' \in \mathbb{R}^p$  at time  $t$ :

$$\mathbf{y}_t = \mathbf{B}\mathbf{f}_t + \mathbf{u}_t, \quad t = 1, \dots, T, \quad (1)$$

where  $\mathbf{f}_t \in \mathbb{R}^K$  is a  $K$ -dimensional vector of common factors,  $\mathbf{B} = (\mathbf{b}'_1, \dots, \mathbf{b}'_p)' \in \mathbb{R}^{p \times K}$  is a factor loading matrix with  $\mathbf{b}_j \in \mathbb{R}^K$  being the factor loading of the  $j$ th variable on the latent factor  $\mathbf{f}_t$ , and  $\mathbf{u}_t$  is an idiosyncratic error vector. In the above model,  $\mathbf{y}_t$  is the only observable variable, while  $\mathbf{B}$  is a matrix of unknown parameters, and  $(\mathbf{f}_t, \mathbf{u}_t)$  are latent random variables. Without loss of generality, we assume  $E(\mathbf{f}_t) = E(\mathbf{u}_t) = \mathbf{0}$  and  $\mathbf{f}_t$  and  $\mathbf{u}_t$  are uncorrelated. Then, the model implied covariance structure is

$$\boldsymbol{\Sigma} = \mathbf{B} \text{cov}(\mathbf{f}_t) \mathbf{B}' + \boldsymbol{\Sigma}_u,$$

where  $\boldsymbol{\Sigma} = E(\mathbf{y}_t \mathbf{y}_t')$  and  $\boldsymbol{\Sigma}_u = E(\mathbf{u}_t \mathbf{u}_t')$ . Observe that  $\mathbf{B}$  and  $\mathbf{f}_t$  are not individually identifiable, since  $\mathbf{B} \mathbf{f}_t = \mathbf{B} \mathbf{H} \mathbf{H}' \mathbf{f}_t$  for any orthogonal matrix  $\mathbf{H}$ . To this end, an identifiability condition is imposed:

$$\text{cov}(\mathbf{f}_t) = \mathbf{I}_K \text{ and } \mathbf{B}' \boldsymbol{\Sigma}_u^{-1} \mathbf{B} \text{ is diagonal,} \quad (2)$$

which is a common assumption in the literature (Bai and Li 2012; Bai and Liao 2013).

Assume that we are only interested in a subset  $S$  among a total of  $p$  variables in model (1). We aim to obtain an efficient estimator of

$$\boldsymbol{\Sigma}_S = \mathbf{B}_S \mathbf{B}_S' + \boldsymbol{\Sigma}_{u,S},$$

the covariance matrix of the  $s$  variables in  $S$ , where  $\mathbf{B}_S$  is the sub-matrix of  $\mathbf{B}$  with row indices in  $S$  and  $\boldsymbol{\Sigma}_{u,S}$  is the submatrix of  $\boldsymbol{\Sigma}_u$  with row and column indices in  $S$ . As mentioned above, the existing literature uses the following conventional method:

- Method 1: Use solely the  $s$  variables in the set  $S$  to estimate common factors  $\mathbf{f}_b$ , the loading matrix  $\mathbf{B}_S$ , the idiosyncratic matrix  $\boldsymbol{\Sigma}_{u,S}$ , and the covariance matrix  $\boldsymbol{\Sigma}_S$ .

This idea is apparently strongly influenced by the nonparametric estimation of the covariance matrix and ignores a large portion of the available data in the other  $p - s$  variables. An intuitively more efficient method is

- Method 2: Use all the  $p$  variables to obtain estimators of  $\mathbf{f}_b$ , the loading matrix  $\mathbf{B}$ , the idiosyncratic matrix  $\boldsymbol{\Sigma}_u$ , and the entire covariance matrix  $\boldsymbol{\Sigma}$ , and then restrict them to the variables of interest. This is the same as estimating  $\mathbf{f}_t$  using all variables, and then estimating  $\mathbf{B}_S$  and  $\boldsymbol{\Sigma}_{u,S}$  based on the model (1) and the subset  $S$  with  $\mathbf{f}_t$  being estimated (observed), and obtaining a plug-in estimator of  $\boldsymbol{\Sigma}_S$ .

We will show that Method 2 is more efficient than Method 1 in the estimation of  $\mathbf{f}_t$  and  $\boldsymbol{\Sigma}_S$  as more auxiliary data information is incorporated. By treating common factor as an unknown parameter, we calculate its Fisher information that grows with more data being used in Method 2. In this case, a more efficient factor estimate can be obtained, for example, through weighted principal component (WPC) method (Bai and Liao 2013). The advantage of factor estimation is further carried over to the estimation of  $\boldsymbol{\Sigma}_S$  by Method 2 in terms of its convergence rate. Moreover, if the number of total variables is sufficiently large, Method 2 is proven to perform as well as an “oracle method,” which observes all latent factors. This lends further support to our aforementioned claim of “blessing of dimensionality.” Such a best possible rate improvement is new to the existing literature, and counted as another

contribution of this article. All these conclusions hold when the number of factors  $K$  is assumed to be fixed and known, while  $s$ ,  $p$ , and  $T$  all tend to infinity.

The idea of using data as much as possible brings computational challenges. Fortunately, we observe that all the  $p$  variables are controlled by the same group of latent factors. Having said that, we can actually *split  $p$  variables* into smaller groups, and then use each group to estimate latent factors. The final factor estimate is obtained by averaging over these repeatedly estimated factors. Obviously, this divide-and-conquer algorithm can be implemented in a parallel computing environment, and thus produces factor estimators in a much more efficient way. On the other hand, our theory illustrates that this new method performs as well as the “pooled analysis,” where we run the method over the whole dataset. Simulation studies further demonstrate the boosted computational speed and satisfactory statistical performance.

The rest of the article is organized as follows. We compare the Fisher information of the factors by the two methods in Section 2. Section 3 describes the WPC method. As a main result, the convergence rates of Different estimators of  $\Sigma_S$  are further compared in Section 4 under various norms. Section 5 introduces the divide-and-conquer method for accelerating computation, while Section 6 presents all simulation results. Section 7 gives a microarray data example to illustrate our proposal. All technical proofs are delegated to the Appendix.

For any vector  $\mathbf{a}$ , let  $\mathbf{a}_S$  denote a sub-vector of  $\mathbf{a}$  with indices in  $S$ . Denote  $\|\mathbf{a}\|$  the Euclidean norm of  $\mathbf{a}$ . For a symmetric matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$ , let  $\mathbf{A}_{I,J}$  be the submatrix of  $\mathbf{A}$  with row and column indices in  $I$  and  $J$ , respectively. We write  $\mathbf{A}_S$  for  $\mathbf{A}_{S,S}$  for simplicity. Let  $\lambda_j(\mathbf{A})$  be the  $j$ th largest eigenvalue of  $\mathbf{A}$ . Denote  $\|\mathbf{A}\| = \max\{|\lambda_1(\mathbf{A})|, |\lambda_d(\mathbf{A})|\}$  the operator norm of  $\mathbf{A}$ ,  $\|\mathbf{A}\|_{\max} = \max_{ij} |a_{ij}|$  the max-norm of  $\mathbf{A}$ , where  $a_{ij}$  is the  $(i, j)$ th entry of  $\mathbf{A}$ ,  $\|\mathbf{A}\|_1 = \max_i \sum_{j=1}^d |a_{ij}|$  the  $L_1$  norm of  $\mathbf{A}$ ,  $\|\mathbf{A}\|_F = \sqrt{\text{tr}(\mathbf{A}'\mathbf{A})}$  the Frobenius norm of  $\mathbf{A}$ , and  $\|\mathbf{A}\|_{\mathbf{M}} = d^{-1/2} \|\mathbf{M}^{-1/2} \mathbf{A} \mathbf{M}^{-1/2}\|_F$  the relative norm of  $\mathbf{A}$  to  $\mathbf{M}$ , where the weight matrix  $\mathbf{M}$  is assumed to be positive definite. For a non-square matrix  $\mathbf{C}$ , let  $\mathbf{C}_S$  be the submatrix of  $\mathbf{C}$  with row indices in  $S$ .

## 2. Fisher Information of Common Factor

In this section, we treat the vector of common factors as a fixed unknown parameter, and compute its Fisher information matrices based on Method 1 and Method 2. In the computation, the loading matrix  $\mathbf{B}$  is treated as deterministic in Proposition 2. In Proposition 3, the Fisher information is computed for each given  $\mathbf{B}$  and then averaged over  $\mathbf{B}$  by regarding it as a realization of a chance process, which bypasses the block diagonal assumption needed without taking average over  $\mathbf{B}$ . In other sections, we adopt the convention regarding the factors as random and  $\mathbf{B}$  as fixed. We start by calculating the Fisher information of  $\boldsymbol{\theta}_t := \mathbf{B}\mathbf{f}_t$ , which serves as an intermediate step in obtaining that for  $\mathbf{f}_t$ . For simplicity of notation, time  $t$  is suppressed in  $(\mathbf{y}_t, \mathbf{f}_t, \mathbf{u}_t, \boldsymbol{\theta}_t)$  so that it becomes  $(\mathbf{y}, \mathbf{f}, \mathbf{u}, \boldsymbol{\theta})$  in this section.

Given a general density function of  $\mathbf{y}$ , denoted as  $h(\mathbf{y}; \boldsymbol{\theta})$ , the Fisher information of  $\boldsymbol{\theta}$  contained in full data is given by

$$I_p(\boldsymbol{\theta}) = E\left[\left(\frac{\partial \log h(\mathbf{y}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right)\left(\frac{\partial \log h(\mathbf{y}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right)'\right].$$

When only data in  $S$  is used, the Fisher information of  $\boldsymbol{\theta}_S$  is given by

$$I_S(\boldsymbol{\theta}_S) = E\left[\left(\frac{\partial \log h_S(\mathbf{y}_S; \boldsymbol{\theta}_S)}{\partial \boldsymbol{\theta}_S}\right)\left(\frac{\partial \log h_S(\mathbf{y}_S; \boldsymbol{\theta}_S)}{\partial \boldsymbol{\theta}_S}\right)'\right].$$

where  $h_S$  is the marginal density of  $\mathbf{y}_S$  for the target set of variable  $S$ . Our first proposition shows that  $\{I_p(\boldsymbol{\theta})\}_S$ , the submatrix of  $I_p(\boldsymbol{\theta})$  restricted on  $S$ , dominates  $I_S(\boldsymbol{\theta}_S)$  under a mild condition.

*Proposition 1.* If  $h(\mathbf{y}; \boldsymbol{\theta}) = h(\mathbf{y} - \boldsymbol{\theta})$  and the density function  $h(\mathbf{y} - \boldsymbol{\theta})$  satisfies the following regularity condition:

$$\nabla_{\mathbf{y}_S} \int h(\mathbf{y}_S - \boldsymbol{\theta}_S, \mathbf{y}_{S^c} - \boldsymbol{\theta}_{S^c}) d\mathbf{y}_{S^c} = \int \nabla_{\mathbf{y}_S} h(\mathbf{y}_S - \boldsymbol{\theta}_S, \mathbf{y}_{S^c} - \boldsymbol{\theta}_{S^c}) d\mathbf{y}_{S^c}, \quad (3)$$

then  $\{I_p(\boldsymbol{\theta})\}_S \succeq I_S(\boldsymbol{\theta}_S)$  in the sense that  $\{I_p(\boldsymbol{\theta})\}_S - I_S(\boldsymbol{\theta}_S)$  is positive semidefinite.

The regularity condition (3) is fairly mild, as illustrated in the following examples.

*Example 1.* In model 1, if  $\mathbf{u}_S$  and  $\mathbf{u}_{S^c}$  are independent, then (3) holds.

*Example 2.* If  $\mathbf{y}$  follows an elliptical distribution that

$$h(\mathbf{y}; \boldsymbol{\theta}) \propto g((\mathbf{y} - \boldsymbol{\theta})' \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\theta})),$$

where the mapping function  $g(t) : [0, \infty) \rightarrow [0, \infty)$  satisfies that  $|g'(t)| \leq cg(t)$  for some positive constant  $c$ , and  $E|\mathbf{y}| < \infty$ , then (3) holds. Example 2 includes some commonly used multivariate distributions as its special cases, for example, the multivariate normal distribution and the multivariate  $t$ -distribution with degrees of freedom greater than 1. The proof is given in Appendix A.2.

We next compute the Fisher information of  $\mathbf{f}$  based on the full dataset, denoted as  $I(\mathbf{f})$ , and the partial dataset restricted on  $S$ , denoted as  $I_S(\mathbf{f})$ . This can be done easily by noting that  $I(\mathbf{f}) = \mathbf{B}' I_p(\boldsymbol{\theta}) \mathbf{B}$ . Indeed, the WPC estimators used in Methods 1 and 2 achieve such efficiency since their asymptotic variances are proven to be the inverse of  $I(\mathbf{f})$  and  $I_S(\mathbf{f})$ , respectively; see Remark 1.

Proposition 2 shows that  $I(\mathbf{f})$  dominates  $I_S(\mathbf{f})$ , if  $I_p(\boldsymbol{\theta})$  is block-diagonal, that is,  $\{I_p(\boldsymbol{\theta})\}_{S, S^c} = \mathbf{0}$ . Hence, common factors can be estimated more efficiently using additional data  $\mathbf{y}_{S^c}$ . The above block-diagonal condition implies that the idiosyncratic error of additional variables

cannot be confounded with that of the variables-of-interest. For example, if  $\mathbf{u}$  is normal, then  $\{I_p(\boldsymbol{\theta})\}_{S,S^c} = \mathbf{0}$  indeed requires that  $\mathbf{u}_S$  is independent of  $\mathbf{u}_{S^c}$ .

*Proposition 2.* Under condition (3), if  $\{I_p(\boldsymbol{\theta})\}_{S,S^c} = \mathbf{0}$ ,  $I(\mathbf{f}) = I_S(\mathbf{f})$ .

So far we treat  $\mathbf{B}$  as being deterministic. Rather, Proposition 3 regards  $\{\mathbf{b}_j\}$  as a realization of a chance process. Under this assumption, the expectation of  $I(\mathbf{f})$  over  $\mathbf{B}$  is shown to always dominate that of  $I_S(\mathbf{f})$ . In other words, we can claim that averaging over loading matrices, a larger dataset contains more information about the unknown factors.

*Proposition 3.* If  $\{\mathbf{B}_i\}_{i=1}^P$  are iid random loadings with  $E(\mathbf{b}_j) = \mathbf{0}$  and (3) holds, then  $E[I(\mathbf{f})] \geq E[I_S(\mathbf{f})]$ , where the expectation is taken with respect to the distribution of  $\mathbf{B}$ .

### 3. Efficient Estimation of Common Factor

In this section, we construct an efficient estimator of the common factors by showing that its asymptotic variance is exactly the inverse of its Fisher information. This together with the arguments in Section 2 enables us to draw a conclusion that using more data results in a more efficient factor estimator with a smaller asymptotic variance.

From a least-square perspective, when the loading matrix  $\mathbf{B}$  is known,  $\mathbf{f}_t$  can be estimated by the weighted least-squares:  $\operatorname{argmin}_{\mathbf{f}_t \in \mathbb{R}^K} \sum_{t=1}^T (\mathbf{y}_t - \mathbf{B}\mathbf{f}_t)' \boldsymbol{\Sigma}_u^{-1} (\mathbf{y}_t - \mathbf{B}\mathbf{f}_t)$ . In the high-dimensional setting ( $p \gg T$ ), we assume  $\boldsymbol{\Sigma}_u$  is a sparse matrix and define its sparsity measurement as

$$m_p = \max_{i \leq p} \sum_{j \neq i} I(\sigma_{u,ij} \neq 0), \text{ where } \sigma_{u,ij} \text{ is the } (i, j) \text{th entry of } \boldsymbol{\Sigma}_u. \quad (4)$$

In particular, we assume the following sparsity condition:

$$m_p = o\left(\min\left\{\frac{1}{p^{1/4}} \sqrt{\frac{T}{\log p}}, p^{1/4}\right\}\right) \text{ and } \sum_{i=1}^p \sum_{j \neq i} I(\sigma_{u,ij} \neq 0) = O(p). \quad (5)$$

Now, we propose to solve the following constrained weighted least-square problem:

$$\begin{aligned} (\hat{\mathbf{B}}, \hat{\mathbf{f}}_1, \dots, \hat{\mathbf{f}}_T) &= \operatorname{argmin}_{\mathbf{B}, \mathbf{f}_t} \sum_{t=1}^T (\mathbf{y}_t - \mathbf{B}\mathbf{f}_t)' \tilde{\boldsymbol{\Sigma}}_u^{-1} (\mathbf{y}_t - \mathbf{B}\mathbf{f}_t), \\ \text{subject to } &\frac{1}{T} \sum_{t=1}^T \mathbf{f}_t \mathbf{f}_t' = \mathbf{I}_K; \mathbf{B}' \tilde{\boldsymbol{\Sigma}}_u^{-1} \mathbf{B} \text{ is diagonal,} \end{aligned} \quad (6)$$

where  $\tilde{\Sigma}_u$  is a regularized estimator of  $\Sigma_u$  to be discussed later. The above constraint is a sample analog of the identifiability condition (2). The involvement of the weight  $\tilde{\Sigma}_u^{-1}$  is to account for the heterogeneity among the data and leads to more efficient estimation of  $(\mathbf{B}, \mathbf{f}_t)$  (Choi 2012; Bai and Liao 2013).

Indeed, an initial estimator  $\tilde{\Sigma}_u$  of the idiosyncratic matrix  $\Sigma_u$  is needed for solving the constrained weighted least-square problem. We propose to obtain such an estimator by the following procedure, which is in the same spirit as the estimation of the idiosyncratic matrix in the POET method (Fan, Liao, and Mincheva 2013). Let  $\mathbf{S}_y = T^{-1} \sum_{t=1}^T (\mathbf{y}_t - \bar{\mathbf{y}})(\mathbf{y}_t - \bar{\mathbf{y}})'$  be the sample covariance of  $\mathbf{y}$  and  $\{(\lambda_i, \zeta_i)\}_{i=1}^p$  be eigen-pairs of  $\mathbf{S}_y$  with  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$ . Denote  $\mathbf{R} = \mathbf{S}_y - \sum_{i=1}^K \lambda_i \zeta_i \zeta_i'$ . We estimate  $\Sigma_u$  by  $\hat{\Sigma}_u$  whose  $(i, j)$ th entry

$$\hat{\sigma}_{u,ij} = \begin{cases} r_{ij}, & \text{for } i = j, \\ s_{ij}(r_{ij}), & \text{for } i \neq j, \end{cases} \text{ where } \mathbf{R} = (r_{ij}),$$

$s_{ij}(r_{ij})$  is a general entry-wise thresholding function (Antoniadis and Fan 2001) such that  $s_{ij}(z) = 0$  if  $|z| \leq \tau_{ij}$  and  $|s_{ij}(z) - z| = \tau_{ij}$  for  $|z| > \tau_{ij}$ . In our article, we choose hard-thresholding even though SCAD (Fan and Li 2001) and MCP (Zhang 2010) are also applicable. We specify the entry-wise thresholding level as

$$\tau_{ij}(p) = C\sqrt{r_{ii}r_{jj}}\omega(p), \text{ where } \omega(p) = \sqrt{\frac{\log p}{T}} + \frac{1}{\sqrt{p}}, \quad (7)$$

and  $C$  is a constant chosen by cross-validation. The thresholding parameter  $C_\omega(p)$  is applied to the correlation matrix. This is similar to the adaptive thresholding estimator for a general covariance matrix (Rothman, Levina, and Zhu 2009), where the entry-wise thresholding level depends on  $p$ .

With  $\tilde{\Sigma}_u$  being the thresholding estimator described above, the constrained weighted least-square problem (6) can be solved by the weighted principal component (WPC) method. The solution is given by

$$\hat{\mathbf{F}} = (\hat{\mathbf{f}}_1, \dots, \hat{\mathbf{f}}_T)' \text{ and } \hat{\mathbf{B}}' = T^{-1} \mathbf{Y} \hat{\mathbf{F}}, \quad (8)$$

where  $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_T)$  and the columns of  $\hat{\mathbf{F}}$  are the eigenvectors corresponding to the largest  $K$  eigenvalues of the  $T \times T$  matrix  $\sqrt{T} \mathbf{Y}' \tilde{\Sigma}_u^{-1} \mathbf{y}$  (Bai and Liao 2013).

In the following, we give a result showing that the WPC estimator is asymptotically efficient. Indeed, Bai and Liao (2013) derive the asymptotic normality of  $\hat{\mathbf{f}}_t$  under the following conditions:

- i. All eigenvalues of  $\mathbf{B}'\mathbf{B}/p$  are bounded away from zero and infinity as  $p \rightarrow \infty$ ;
- ii. There exists a  $K \times K$  diagonal matrix  $\mathbf{Q}$  such that  $\mathbf{B}'\boldsymbol{\Sigma}_u^{-1}\mathbf{B}/p \rightarrow \mathbf{Q}$ . In addition, the diagonal elements of  $\mathbf{Q}$  are distinct and bounded away from infinity.
- iii. For each fixed  $t \leq T$ ,  $(\mathbf{B}'\boldsymbol{\Sigma}_u^{-1}\mathbf{B})^{-1/2}\mathbf{B}'\boldsymbol{\Sigma}_u^{-1}\mathbf{u}_t \xrightarrow{d} N(\mathbf{0}, \mathbf{I}_K)$ , as  $p \rightarrow \infty$ , together with the sparsity assumption (5), and some additional regularity conditions given in Section A.1. When  $p \log p = o(T)$ , it is shown that

$$\sqrt{p}(\hat{\mathbf{f}}_t - \mathbf{H}\mathbf{f}_t) \xrightarrow{D} N(\mathbf{0}, \mathbf{Q}^{-1}), \quad (9)$$

where  $\mathbf{H}$  is a specific rotation matrix given by

$$\mathbf{H} = \hat{\mathbf{V}}^{-1}\hat{\mathbf{F}}'\mathbf{F}\mathbf{B}'\tilde{\boldsymbol{\Sigma}}_u^{-1}\mathbf{B}/T, \quad (10)$$

and  $\hat{\mathbf{V}}$  is a  $K \times K$  diagonal matrix of the largest  $K$  eigenvalues of  $\mathbf{Y}'\tilde{\boldsymbol{\Sigma}}_u^{-1}\mathbf{y}/T$ . The rotation matrix  $\mathbf{H}$  is introduced here so that  $\mathbf{H}\mathbf{f}_t$  is an identifiable quantity from the data. See more discussion about the identifiability in Remark 2.

Condition (i) is a ‘‘pervasive condition’’ requiring that the common factors affect a nonnegligible fraction of subjects. This is a common assumption for the principal components based methods (Fan, Liao, and Mincheva 2011; Bai and Liao 2013). In condition (ii),  $\mathbf{B}'\boldsymbol{\Sigma}_u^{-1}\mathbf{B}$  is indeed the Fisher information (under Gaussian errors) contained in  $p$  variables, while the limit  $\mathbf{Q}$  can be viewed as an average information for each variable. Hence, the asymptotic normality in (9) shows that  $\hat{\mathbf{f}}_t$  is efficient as its asymptotic variance attains the inverse of the (averaged) Fisher information.

*Remark 1.* The results in Section 2 together with (9) imply that Method 2 is in general better than Method 1 in the estimation of common factors. To explain why, we consider two different cases here. When  $p$  is an order of magnitude larger than  $s$ , where  $s$  is the number of variables of interest. Method 2 produces a better estimator of factors with a faster convergence rate. Even when  $p$  and  $s$  diverge at the same speed, the factor estimator based on Method 2 is shown to possess a smaller asymptotic variance, as long as  $\boldsymbol{\Sigma}_{u,S,S^c} = \mathbf{0}$ . Recall that  $\mathbf{B}'\boldsymbol{\Sigma}_u^{-1}\mathbf{B} = \mathbf{I}(\mathbf{f})$  and  $\mathbf{B}'_S\boldsymbol{\Sigma}_{u,S}^{-1}\mathbf{B}_S = \mathbf{I}_S(\mathbf{f})$  under Gaussian errors, and they also correspond to the inverse of the asymptotic variance given by Methods 1 and 2, respectively. Then, Proposition 2 implies that Method 2 has a smaller asymptotic variance, if  $\boldsymbol{\Sigma}_{u,S,S^c} = \mathbf{0}$ . Alternatively, if  $\mathbf{B}$  is treated as being random, Proposition 3 immediately implies that



$E(\mathbf{B}'_S \boldsymbol{\Sigma}_{u,S}^{-1} \mathbf{B}_S) \geq E(\mathbf{B}' \boldsymbol{\Sigma}^{-1} \mathbf{B})$ . Therefore, even without the block diagonal assumption, Method 2 produces a more efficient factor estimate on average.

#### 4. Covariance Matrix Estimation

One primary goal in this article is to obtain an accurate estimator of the covariance matrix  $\boldsymbol{\Sigma}_S = E(\mathbf{y}_S \mathbf{y}'_S)$  for the variables-of-interest. In this section, we compare three Different estimation methods, namely, Methods 1, 2, and Oracle Method, in terms of their rates of convergence (under various norms). Obviously, these rates depend on how accurately the realized factors are estimated as demonstrated later.

Below we describe these three methods in full details.

- Method 1:
  - i. Use solely the data in the subset  $S$  to obtain estimators of the realized factors  $\hat{\mathbf{F}}^{(1)}$  and the loading matrix  $\hat{\mathbf{B}}_1 = T^{-1} \mathbf{Y}_S \hat{\mathbf{F}}^{(1)}$  based on (8);
  - ii. Let  $(\hat{\mathbf{f}}_t^{(1)})'$  be the  $t$ th row of  $\hat{\mathbf{F}}^{(1)}$ ,  $(\hat{\mathbf{b}}_i^{(1)})'$  be the  $i$ th row of  $\hat{\mathbf{B}}_1$ ,  $\hat{u}_{it} = y_{it} - (\hat{\mathbf{b}}_i^{(1)})' \hat{\mathbf{f}}_t^{(1)}$ , and  $\hat{\sigma}_{ij} = \frac{1}{T} \sum_{t=1}^T \hat{u}_{it} \hat{u}_{jt}$ . The  $(i, j)$ th entry of the idiosyncratic matrix estimator  $\hat{\boldsymbol{\Sigma}}_{u,S}^{(1)}$  of  $\boldsymbol{\Sigma}_{u,S}$  is given by thresholding  $\hat{\sigma}_{ij}$  at the level of  $C \hat{\theta}_{ij}^{1/2} \omega(s)$ , where  $\omega(s)$  is defined in (7) and  $\hat{\theta}_{ij} = \frac{1}{T} \sum_{t=1}^T (\hat{u}_{it} \hat{u}_{jt} - \hat{\sigma}_{ij})^2$ ;
  - iii. The final estimator is given by  $\hat{\boldsymbol{\Sigma}}_S^{(1)} = \hat{\mathbf{B}}_1 \hat{\mathbf{B}}_1' + \hat{\boldsymbol{\Sigma}}_{u,S}^{(1)}$ .
- Method 2:
  - i. Use all  $p$  variables to obtain the estimate  $\hat{\mathbf{F}}^{(2)}$  as given in (8) for the realized factors and then estimate the loading  $\mathbf{B}_S$  by  $\hat{\mathbf{B}}_2 = T^{-1} \mathbf{Y}_S \hat{\mathbf{F}}^{(2)}$ ;
  - ii. Follow the same procedure as in Method 1 to obtain the estimator  $\hat{\boldsymbol{\Sigma}}_{u,S}^{(2)}$  but based on  $\hat{\mathbf{F}}^{(2)}$  and  $\hat{\mathbf{B}}_2$ ;
  - iii. The final estimator is given by  $\hat{\boldsymbol{\Sigma}}_S^{(2)} = \hat{\mathbf{B}}_2 \hat{\mathbf{B}}_2' + \hat{\boldsymbol{\Sigma}}_{u,S}^{(2)}$ .
- Oracle Method:
  - i. Estimate the loading by  $\hat{\mathbf{B}}_o = T^{-1} \mathbf{Y}_S \mathbf{F}$ , where  $\mathbf{F} = (\mathbf{f}_1, \dots, \mathbf{f}_T)'$  are the true factors.
  - ii. The idiosyncratic matrix estimator  $\hat{\boldsymbol{\Sigma}}_{u,S}^o$  is given by the same procedure as in Method 1, with  $\hat{\mathbf{b}}_i^{(1)}$  and  $\hat{\mathbf{f}}_t^{(1)}$  being replaced by  $\hat{\mathbf{b}}_i^o$  and  $\mathbf{f}_t$  respectively.
  - iii. The final estimator is given by  $\hat{\boldsymbol{\Sigma}}_S^o = \hat{\mathbf{B}}_o \hat{\mathbf{B}}_o' + \hat{\boldsymbol{\Sigma}}_{u,S}^o$ .

Theorem 1 depicts the estimation accuracy of  $\Sigma_S$  by the above three methods with respect to the following measurements:

$$\|\widehat{\Sigma}_S - \Sigma_S\|_{\Sigma_S}, \|\widehat{\Sigma}_S - \Sigma_S\|_{\max}, \|\widehat{\Sigma}_S^{-1} - \Sigma_S^{-1}\|,$$

where  $\|\widehat{\Sigma}_S - \Sigma_S\|_{\Sigma_S} = p^{-1/2} \|\Sigma_S^{-1/2} \widehat{\Sigma}_S \Sigma_S^{-1/2} - \mathbf{I}_S\|_F$  is a norm of the relative errors. Note that

the results of Fan, Liao, and Mincheva (2013) cannot be directly used here since we employ the *weighted* principal component analysis to estimate the unobserved factors. This is expected to be more accurate than the ordinary principal component analysis, as shown in Bai and Liao (2013). Indeed, the technical proofs for our results are technically more involved than those by Fan, Liao, and Mincheva (2013).

We assume that  $s$  is much less than  $p$ , that is,  $s = o(p)$ , but both tend to infinity. Under the pervasive condition (i),  $\|\Sigma_S\| \asymp cs$  and therefore diverges. For this reason, we consider the relative norm  $\|\widehat{\Sigma}_S - \Sigma_S\|_{\Sigma_S}$ , instead of  $\|\widehat{\Sigma}_S - \Sigma_S\|$ , and the operator norm  $\|\widehat{\Sigma}_S^{-1} - \Sigma_S^{-1}\|$  for estimating the inverse. In addition, we consider another element-wise max norm  $\|\widehat{\Sigma}_S - \Sigma_S\|_{\max}$ . We show that if  $p$  is large with respect to  $s$  and  $T$ , Method 2 performs as well as the Oracle Method, both of which outperform Method 1. As a consequence, even if we are only interested in the covariance matrix of a small subset of variables, we should use all the data to estimate the common factors, which ultimately improves the estimation of  $\Sigma_S$ . In particular, we are able to specify an explicit regime of  $(s, p)$  under which the improvements are substantial. However, when  $s \asymp p$ , that is, they are in the same order, using more data does not show as dramatic improvements for estimating  $\Sigma_S$ . This is expected and will be clearly seen in the simulation section.

Before stating Theorem 1, we need a few preliminary results: Lemmas 1–3. Specifically, Lemma 1 presents the uniform convergence rates of the factor estimates by Methods 1 and 2. Based on that, Lemmas 2 and 3 further derive the estimation accuracy of factor loadings and idiosyncratic matrix by the three methods, respectively. These results together lead to the estimation error rates of  $\Sigma_S$  in Theorem 1 w.r.t. three measures defined above. Additional Lemmas supporting the proof are given in the Appendix. Again, these kinds of results cannot be obtained directly from Fan, Liao, and Mincheva (2013) due to our use of WPC.

*Lemma 1.* Suppose that conditions (i), (ii), the sparsity condition (5), and additional regularity conditions (iv)–(vii) in Section A.1 hold for both  $s$  and  $p$ . If  $p \log p = o(T)$  and  $T = o(s^2)$ , then we have

$$\begin{aligned} \max_{t \leq T} \|\widehat{\mathbf{f}}_t^{(1)} - \mathbf{H}_1 \mathbf{f}_t\| &= O_P\left(\frac{1}{\sqrt{T}} + \frac{T^{1/4}}{\sqrt{s}}\right) \text{ and} \\ \max_{t \leq T} \|\widehat{\mathbf{f}}_t^{(2)} - \mathbf{H}_2 \mathbf{f}_t\| &= O_P\left(\frac{1}{\sqrt{T}} + \frac{T^{1/4}}{\sqrt{p}}\right). \end{aligned}$$

where  $\mathbf{H}_1 = \widehat{\mathbf{V}}_1^{-1} \widehat{\mathbf{F}}^{(1)'} \mathbf{F} \mathbf{B}'_S \widetilde{\Sigma}_{u,S}^{-1} \mathbf{B}_S / T$ ,  $\mathbf{H}_2 = \widehat{\mathbf{V}}_2^{-1} \widehat{\mathbf{F}}^{(2)'} \mathbf{F} \mathbf{B}'_S \widetilde{\Sigma}_{u,S}^{-1} \mathbf{B}_S / T$ ,  $\widehat{\mathbf{V}}_1$  is the diagonal matrix of the largest  $K$  eigenvalues of  $\mathbf{Y}'_S \widetilde{\Sigma}_{u,S}^{-1} \mathbf{Y}_S / T$  and  $\widehat{\mathbf{V}}_2$  is the diagonal matrix of the largest  $K$  eigenvalues of  $\mathbf{Y}' \widetilde{\Sigma}_u^{-1} \mathbf{y} / T$ .

*Remark 2.*  $\mathbf{H}_1$  and  $\mathbf{H}_2$  correspond to the rotation matrix  $\mathbf{H}$  defined in (10) using Methods 1 and 2, respectively. Recall that  $\mathbf{F} = (\mathbf{f}_1, \dots, \mathbf{f}_T)'$ , then  $\mathbf{H}\mathbf{f}_t = T^{-1} \widehat{\mathbf{V}}^{-1} \widehat{\mathbf{F}} (\mathbf{B}\mathbf{f}_1, \dots, \mathbf{B}\mathbf{f}_T)' \widetilde{\Sigma}_u^{-1} \mathbf{B}\mathbf{f}_t$ . Note that  $\mathbf{H}\mathbf{f}_t$  only depends on quantities  $\mathbf{V}^{-1} \widehat{\mathbf{F}}$ ,  $\widetilde{\Sigma}_u^{-1}$  and the *identifiable* component  $\{\mathbf{B}\mathbf{f}_t\}_{t=1}^T$ . Therefore, there is no identifiability issue regarding  $\mathbf{H}\mathbf{f}_t$ . In other words, even though  $\mathbf{f}_t$  itself may not be identifiable, an identifiable rotation of  $\mathbf{f}_t$  can be consistently estimated by  $\widehat{\mathbf{f}}_t$ .

Lemma 1 implies that Method 2 produces a better factor estimate if

$$0.5 < \gamma_s < 1.5 \leq \gamma_p < 2,$$

by representing  $s$  and  $p$  as  $s \asymp T^\gamma$  and  $p \asymp T^\gamma$ .

It is not surprising that the estimation accuracy of loading matrix also varies among these three methods as shown in Lemma 2.

*Lemma 2.* Under conditions of Lemma 1,

$$\begin{aligned} \max_{i \leq s} \|\widehat{\mathbf{b}}_i^{(1)} - \mathbf{H}_1 \mathbf{b}_i\| &= O_P(w_1), \\ \text{where } w_1 &:= \frac{1}{\sqrt{s}} + \sqrt{\frac{\log s}{T}}, \max_{i \leq s} \|\widehat{\mathbf{b}}_i^{(2)} - \mathbf{H}_2 \mathbf{b}_i\| = O_P(w_2), \\ \text{where } w_2 &:= \frac{1}{\sqrt{p}} + \sqrt{\frac{\log s}{T}}, \max_{i \leq s} \|\widehat{\mathbf{b}}_i^{(2)} - \mathbf{b}_i\| = O_P(w_o), \\ \text{where } w_o &:= \sqrt{\frac{\log s}{T}}. \end{aligned}$$

Similarly, Lemma 2 indicates that Method 2 performs as well as the Oracle Method, both of which are better than Method 1, that is,  $w_2 = w_o < w_1$ , if

$$0.5 < \gamma_s < 1 \leq \gamma_p < 2,$$

by representing  $s$  and  $p$  in the order of  $T$  as above. We remark that the extra terms  $1/\sqrt{s}$  and  $1/\sqrt{p}$  in  $w_1$  and  $w_2$  (in comparison with the oracle rate  $w_o$ ) are due to the factor estimation. Another preliminary result regarding the estimation of the identifiable component  $\mathbf{b}_i \mathbf{f}_i'$  is given in Lemma A.1.

Similar insights can be delivered from Lemma 3 on the estimation of  $\Sigma_{u,S}$ .

*Lemma 3.* Under conditions of Lemma 1, it holds that

$$\begin{aligned} \|\widehat{\Sigma}_{u,S}^{(1)} - \Sigma_{u,S}\| &= O_p(m_s w_1) = \left\| (\widehat{\Sigma}_{u,S}^{(1)})^{-1} - \Sigma_{u,S}^{-1} \right\|, \\ \|\widehat{\Sigma}_{u,S}^{(2)} - \Sigma_{u,S}\| &= O_p(m_s w_2) = \left\| (\widehat{\Sigma}_{u,S}^{(2)})^{-1} - \Sigma_{u,S}^{-1} \right\|, \\ \|\widehat{\Sigma}_{u,S}^o - \Sigma_{u,S}\| &= O_p(m_s w_o) = \left\| (\widehat{\Sigma}_{u,S}^o)^{-1} - \Sigma_{u,S}^{-1} \right\|, \end{aligned}$$

where  $m_s$  is defined as in (4) with  $p$  being replaced by  $s$ .

Now, we are ready to state our main result on the estimation of  $\Sigma_S$  based on the above preliminary results. From Theorem 1, it is easily seen that the comparison of the estimation accuracy of  $\Sigma_S$  among three methods is solely determined by the relative magnitude of  $w_o$ ,  $w_1$ , and  $w_2$ . Therefore, we should use additional variables to estimate the factors if  $p$  is much larger than  $s$  in the sense that  $T \log s = O(p)$  and  $s \log s = o(T)$  (implying  $w_2 = w_o < w_1$ ).

*Theorem 1.* Under conditions of Lemma 1, it holds that

1. For the relative norm,  $\|\widehat{\Sigma}_S^{(1)} - \Sigma_S\|_{\Sigma_S} = O_p(\sqrt{s}w_1^2 + m_s w_1)$ ,  
 $\|\widehat{\Sigma}_S^{(2)} - \Sigma_S\|_{\Sigma_S} = O_p(\sqrt{s}w_2^2 + m_s w_2)$ , and  $\|\widehat{\Sigma}_S^o - \Sigma_S\|_{\Sigma_S} = O_p(\sqrt{s}w_o^2 + m_s w_o)$ .
2. For the max-norm,  $\|\widehat{\Sigma}_S^{(1)} - \Sigma_S\|_{\max} = O_p(w_1)$ ,  $\|\widehat{\Sigma}_S^{(2)} - \Sigma_S\|_{\max} = O_p(w_2)$ , and  
 $\|\widehat{\Sigma}_S^o - \Sigma_S\|_{\max} = O_p(w_o)$ .
3. For the operator norm of the inverse matrix,  $\|(\widehat{\Sigma}_S^{(1)})^{-1} - \Sigma_S^{-1}\| = O_p(m_s w_1)$ ,  
 $\|(\widehat{\Sigma}_S^{(2)})^{-1} - \Sigma_S^{-1}\| = O_p(m_s w_2)$  and  $\|(\widehat{\Sigma}_S^o)^{-1} - \Sigma_S^{-1}\| = O_p(m_s w_o)$ .

*Remark 3.* So far, we assumed that the number of factors  $K$  is fixed and known. A data-driven choice of  $K$  has been extensively studied in the econometrics literature, for example, by Bai and Ng (2002), Kapetanios (2010). To estimate  $K$ , we can adopt the method by Bai and Ng (2002) and propose a consistent estimator of  $K$  (by allowing  $p, T \rightarrow \infty$ ) as follows:

$$\widehat{K} = \underset{0 \leq k \leq N}{\operatorname{argmin}} \log \left\{ \frac{1}{pT} \|\mathbf{Y} - T^{-1} \mathbf{Y} \widehat{\mathbf{F}}_k \widehat{\mathbf{F}}_k'\|_F^2 \right\} + kg(p, T),$$

where  $N$  is a predefined upper bound,  $\widehat{\mathbf{F}}_k$  is a  $T \times k$  matrix whose columns are  $T$  times the eigenvectors corresponding to the largest  $k$  eigenvalues of  $\mathbf{Y}'\mathbf{Y}$ , and  $g(p, T)$  is a penalty function. Two examples suggested by Bai and Ng (2002) are

$$g(T, p) = \frac{p+T}{pT} \log \left( \frac{pT}{p+T} \right) \quad \text{or}$$

$$g(T, p) = \frac{p+T}{pT} \log (\min \{p, T\}).$$

Under our assumptions (i)–(x), all conditions required by theorem 2 of Bai and Ng (2002) hold. Hence, their theorem implies that  $P(\hat{K} = K) \rightarrow 1$ . Then, conditioning on the event that  $\{\hat{K} = K\}$ , our Theorem 1 still holds by replacing  $K$  with  $\hat{K}$ . Other effective methods for selecting the number of factors include the eigen ratio method by Lam and Yao (2012) and Ahn and Horenstein (2013).

*Remark 4.* When  $K$  grows with  $p$  and  $T$ , Fan, Liao, and Mincheva (2013) gave the explicit dependence of the convergence rates on  $K$  for their proposed POET estimator. By adopting their technique, we can obtain the following results:

1.  $\|\widehat{\Sigma}_S^{(1)} - \Sigma_S\|_{\Sigma_S} = O_p(K\sqrt{sw_1^2} + K^3 m_s w_1)$ ,  $\|\widehat{\Sigma}_S^{(2)} - \Sigma_S\|_{\Sigma_S} = O_p(K\sqrt{sw_2^2} + K^3 m_s w_2)$ ,  
 $\|\widehat{\Sigma}_S^o - \Sigma_S\|_{\Sigma_S} = O_p(K\sqrt{sw_o^2} + K^3 m_s w_o)$ ;
2.  $\|\widehat{\Sigma}_S^{(1)} - \Sigma_S\|_{\max} = O_p(K^3 w_1)$ ,  $\|\widehat{\Sigma}_S^{(2)} - \Sigma_S\|_{\max} = O_p(K^3 w_2)$ ,  
 $\|\widehat{\Sigma}_S^o - \Sigma_S\|_{\max} = O_p(K^3 w_o)$ ;
3.  $\|(\widehat{\Sigma}_S^{(1)})^{-1} - \Sigma_S^{-1}\| = O_p(K^3 m_s w_1)$ ,  $\|(\widehat{\Sigma}_S^{(2)})^{-1} - \Sigma_S^{-1}\| = O_p(K^3 m_s w_2)$ ,  
 $\|(\widehat{\Sigma}_S^o)^{-1} - \Sigma_S^{-1}\| = O_p(K^3 m_s w_o)$ .

Again, the rate difference among three types of estimators only depends on  $w_o$ ,  $w_1$ , and  $w_2$ . Therefore, the same conclusion (when  $p$  is much larger than  $s$ , using additional variables improves the estimation of  $\Sigma_S$ ) can still be made even if  $K$  diverges. As long as  $K$  diverges in the rate that  $K = o(\min \{1/(\sqrt{sw_1^2}), 1/(m_s w_1)^{1/3}\})$ ,  $K = o(1/w_1^{1/3})$  or  $K = o(1/(m_s w_1)^{1/3})$ , the same blessing of dimensionality phenomena persist in terms of estimation consistency in relative norm, max norm, or operator norm of the inverse, respectively.

## 5. Divide-and-Conquer Computing Method

As discussed previously, we prefer using auxiliary data information as much as possible even we are only interested in the covariance matrix of some particular set of variables. But this can bring up heavy computational burden. This concern motivates a simple divide-and-conquer scheme that *splits all  $p$  variables in  $\mathbf{Y}$* . Without loss of generality, assume that  $p$  rows of matrix  $\mathbf{Y}$  can be evenly divided into  $M$  groups with  $p/M$  variables in each group. The  $s$  variables of interest can possibly be assigned to Different groups.

Divide-and-Conquer Computation Scheme

1. In the  $m$ th group, obtain the initial estimator  $\tilde{\Sigma}_{u,m}$  by using the adaptive thresholding method as described in Section 3 based on the data in the  $m$ th group only.
2. Denote  $\mathbf{Y}_m$  as the data vector corresponding to the variables in the  $m$ th group and let  $\hat{\mathbf{F}}_m = (\hat{\mathbf{f}}_{m,1}, \dots, \hat{\mathbf{f}}_{m,T})'$ , where its columns are the eigenvectors corresponding to the largest  $K$  eigenvalues of the  $T \times T$  matrix  $\sqrt{T}\mathbf{Y}'_m \tilde{\Sigma}_{u,m}^{-1} \mathbf{Y}_m$ . The computation in the above two steps can be done in a parallel manner.
3. Average  $\{\hat{\mathbf{f}}_{m,t}\}_{m=1}^M$  to obtain a single estimator of  $\mathbf{f}_t$  as

$$\bar{\mathbf{f}}_t = \frac{1}{M} \sum_{m=1}^M \hat{\mathbf{f}}_{m,t}.$$

The loading matrix estimate is given by  $\bar{\mathbf{B}}_S = T^{-1} \mathbf{Y}_S \bar{\mathbf{F}}$ , where  $\bar{\mathbf{F}} = (\bar{\mathbf{f}}_1, \dots, \bar{\mathbf{f}}_T)'$ .

4. The idiosyncratic matrix is estimated as follows. Let  $\bar{\mathbf{f}}'_t$  be the  $t$ th row of  $\bar{\mathbf{F}}$  and  $\bar{\mathbf{b}}'_i$  be the  $i$ th row of  $\bar{\mathbf{B}}_S$ . Let  $\hat{u}_{it} = y_{it} - \bar{\mathbf{b}}'_i \bar{\mathbf{f}}'_t$ ,  $\hat{\sigma}_{ij} = T^{-1} \sum_{t=1}^T \hat{u}_{it} \hat{u}_{jt}$ , and  $\hat{\theta}_{ij} = T^{-1} \sum_{t=1}^T (\hat{u}_{it} \hat{u}_{jt} - \hat{\sigma}_{ij})^2$ . The  $(i, j)$ th entry of  $\Sigma_{u,S}$  is given by thresholding  $\hat{\sigma}_{ij}$  at the level of  $C \hat{\theta}_{ij}^{1/2} \omega(s)$ , where  $\omega(s)$  is defined as in (7) with  $p$  replaced by  $s$ .
5. The final estimator of the covariance matrix is given by

$$\bar{\Sigma}_S = \bar{\mathbf{B}}_S \bar{\mathbf{B}}'_S + \bar{\Sigma}_{u,S}.$$

We show that, if  $M$  is fixed,

$$\begin{aligned} \|\bar{\Sigma}_S - \Sigma_S\|_{\Sigma_S} &= O_p(\sqrt{sw_2^2} + m_s w_2), \\ \|\bar{\Sigma}_S - \Sigma_S\|_{\max} &= O_p(w_2), \\ \|(\bar{\Sigma}_S)^{-1} - \Sigma_S^{-1}\| &= O_p(m_s w_2). \end{aligned}$$

These rates match the rates of  $\hat{\Sigma}_S^{(2)}$  attained by Method 2, where all  $p$  variables are pooled together for the analysis. The proof is given in Appendix A.3. The simulation results in Section 6 further demonstrate that without sacrificing the estimation accuracy, the divide-and-conquer method runs much faster than Method 2. Therefore, the divide-and-conquer method is practically useful when dealing with massive dataset.

The main computational cost of our method comes from taking the inverse of  $\tilde{\Sigma}_{u,m}$ . For our Method 2, where all  $p$  variables are pooled together for the analysis, the computational complexity of the inversion is  $O(p^3)$ . On the other hand, for the divide-and-conquer method,

the corresponding estimator  $\tilde{\Sigma}_{u,m}$  in the  $m$ th group only needs a computational cost of  $\mathcal{O}(p/M^3)$  to be inverted. Then, the total computation complexity is  $\mathcal{O}(p^3/M^2)$ . Hence, the computational speed can be boosted by  $M^2$ -fold. Such a computational acceleration can also be observed from simulation study results in Figure 1(d). Other operations like the eigen-decomposition on the  $T \times T$  matrix  $\sqrt{T}\mathbf{Y}'\tilde{\Sigma}_u^{-1}\mathbf{Y}$  do not have dominating computational cost, as we assume that  $p$  is much larger than  $T$ . When  $M$  grows too fast, the divide-and-conquer method may lose estimation efficiency compared with the pooled analysis (Method 2). However, considering its boost of computation, the divide-and-conquer method is practically useful when dealing with massive dataset.

## 6. Simulations

We use simulated examples to compare the statistical performances of Methods 1, 2, and the Oracle Method. We fix the number of factors  $K = 3$  and repeat 100 simulations for each combination of  $(s, p, T)$ . The loading  $\mathbf{b}_j$ , the factor  $\mathbf{f}_t$  and the idiosyncratic error  $\mathbf{u}_t$  are generated as follows:

- $\{\mathbf{b}_j\}_{j=1}^K$  are iid from  $N_K(\mathbf{0}, 5\mathbf{I}_K)$ .
- $\{\mathbf{f}_t\}_{t=1}^T$  are iid from  $N_K(\mathbf{0}, \mathbf{I}_K)$ .
- $\{\mathbf{u}_t\}_{t=1}^T$  are iid from  $N_p(\mathbf{0}, 50\mathbf{I}_p)$ .

The observations  $\{\mathbf{y}_t\}_{t=1}^T$  are generated from (1) using  $\mathbf{b}_j$ ,  $\mathbf{f}_t$  and  $\mathbf{u}_t$  from the above. Tables 1–4 report the estimation errors of the factors, the loading matrices, and the covariance-of-interest  $\Sigma_S$  in terms of Different measurements.

We see from Tables 1 and 2 that when  $s = 50$  and  $p = 1000, 2000$ , Method 1 performs much worse than Method 2, for both  $T = 200$  and  $T = 400$ . However, when  $s$  increases to 800 with  $p$  being the same, Tables 3 and 4 show that the improvement of Method 2 over Method 1 is less profound. This is expected as the set of interest already contains sufficiently rich information to produce an accurate estimator for realized factors. In general, we note that Method 2 is the most advantageous in the settings where  $s$  is much smaller than  $p$ . In addition, from Tables 3 and 4, we can tell that Method 2 comes closer to the Oracle method as  $p$  grows. In practice, we also observe that the WPC factor estimator performs better than the unweighted PC estimator when  $\mathbf{u}_t$  is heteroscedastic. Due to the space limit, we choose not to present the simulation results in this model.

For further comparison with the divide-and-conquer method, we vary  $T$  from 50 to 500 and set  $(s, p, M)$  as  $s = \lfloor T^{0.6} \rfloor$ ,  $p = \lfloor T^{1.4} \rfloor$ , and  $M = \lfloor T^{0.2} \rfloor$ . Figure 1 shows the estimation errors of the four methods together with the corresponding computational time. Again, when  $p$  is large, Method 2 performs as well as the Oracle Method, both of which greatly outperform Method 1. However, its computation becomes much slower in this case. In contrast, the divide-and-conquer method is much faster, while maintaining comparable

performance as Method 2. In the extreme case that  $p$  is around 6000 ( $T=500$ ), the divide-and-conquer method can boost the speed by nine-fold for Method 2.

## 7. Real Data Example

We use a real data example to illustrate how Different utilization of available variables can affect the inference of the variables of interest. Krug et al. (2012) carried out a gene profiling study among 40 Portuguese and Spanish adults to identify key genetic risk factors for ischemic stroke. Among them, 20 subjects were patients having ischemic stroke and the others were controls. Their gene profiles were obtained using the GeneChip Human Genome U133 Plus 2.0 microarray. The data were available at Gene Expression Omnibus with access name “GSE22255.”

To judge how effectively the gene expression can distinguish ischemic stroke and controls, we applied the Linear Discriminant Analysis (LDA) to this dataset. We randomly chose 10 subjects as the test set and the rest as the training set. We repeated the random splitting for 100 runs. In each run, we selected the set of expressed Differentially (DE) genes with a threshold of over 1.2-fold change and a  $Q$ -value  $< 0.05$ , which is a commonly used quantity to define DE genes (Storey 2002). An LDA rule was then learned from the training set using the selected genes and further applied to the test set for classifying cases and controls. The LDA rule classifies a subject as a case if

$$\hat{\boldsymbol{\delta}}' \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{x} - \bar{\boldsymbol{\mu}}) \geq 0, \quad (11)$$

where  $\hat{\boldsymbol{\delta}} = \hat{\boldsymbol{\mu}}_1 - \hat{\boldsymbol{\mu}}_0 \in \mathbb{R}^s$  is the sample mean difference between the two groups (case-control),  $s$  is the number of selected genes,  $\hat{\boldsymbol{\Sigma}} \in \mathbb{R}^{s \times s}$  is an estimator of the true covariance matrix  $\boldsymbol{\Sigma}$  of the selected genes, and  $\bar{\boldsymbol{\mu}} = (\hat{\boldsymbol{\mu}}_1 + \hat{\boldsymbol{\mu}}_0)/2$ .  $\bar{\boldsymbol{\mu}}$ ,  $\hat{\boldsymbol{\delta}}$ , and  $\hat{\boldsymbol{\Sigma}}$  are obtained from the training set and  $\mathbf{x}$  is the gene expression of subjects in the test set.

As  $s$  can be larger than the sample size, the traditional LDA where  $\hat{\boldsymbol{\Sigma}}$  is the sample covariance is no longer applicable. An alternative method to estimate  $\boldsymbol{\Sigma}$  is adopting the factor model. Factor modeling is widely used in the genomics literature to model the dependencies among genes (Kustra, Shioda, and Zhu 2006; Carvalho et al. 2012). Several factors, like the natural pathway structure (Ogata et al. 2000) can be the latent factors affecting the correlation among genes. A few spiked eigenvalues of the sample covariance in Figure 2 also suggest the existence of potential latent factors in this dataset. Again, there are two ways using the factor model. One way is to use Method 1, where all procedures are done based on the selected genes only. The resulting rule is referred as “LDA-1” in Figure 3. Another way is to use auxiliary data as in Method 2. More specifically, it first uses data from all involved genes and subjects in the training set to estimate the latent factors. These estimated factors are then applied to the set of selected genes, where their loadings and idiosyncratic matrix estimators are obtained. Combining them together produces the covariance matrix estimator, which is still an  $s \times s$  matrix. The resulting rule is referred as



“LDA-2” in Figure 3. Recall that the only difference between the two rules is that they use Different covariance estimators.

Figure 3 plots the average misclassification rates on the test set against the number of factors for the 100 random splits. It is clearly seen that LDA-2 gives better misclassification rates than LDA-1, which is solely due to a Different estimation of the covariance matrix. The results lend further support to our claim that using more data is beneficial.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgments

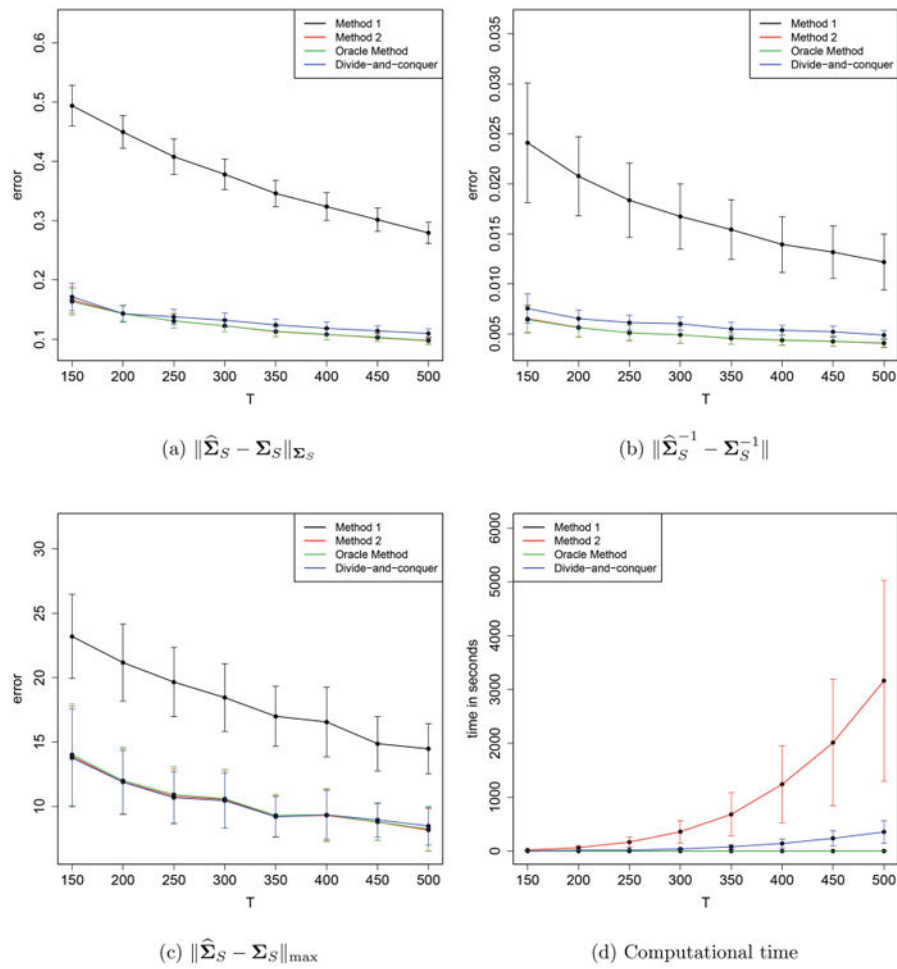
Guang Cheng was on sabbatical at Princeton while this work was carried out and thanks the Princeton ORFE department for its hospitality. The authors thank the editor, the associate editor, and the referees for their constructive comments.

**Funding:** Guang Cheng gratefully acknowledges NSF CAREER Award DMS-1151692, DMS-1418042, Simons Fellowship in Mathematics, Office of Naval Research (ONR N00014-15-1-2331). Jianqing Fan was supported in part by NSF Grants DMS-1206464, DMS-1406266, and NIH grant R01GM100474-4. Quefeng Li was supported by NIH grant 2R01-GM072611-11 as a postdoctoral fellow at Princeton University and supported by NIH grant 2R01-GM047845-25 at UNC-Chapel Hill.

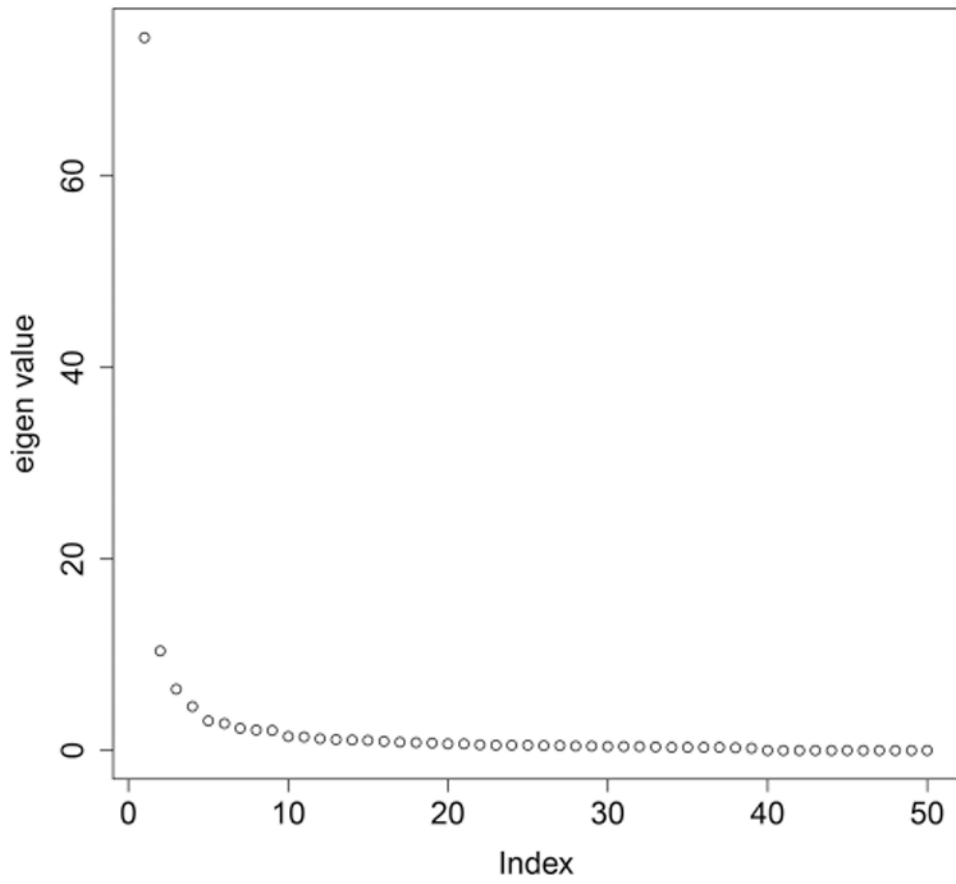
## References

- Ahn SC, Horenstein AR. Eigenvalue Ratio Test for the Number of Factors. *Econometrica*. 2013; 81:1203–1227.
- Antoniadis A, Fan J. Regularized Wavelet Approximations (with discussion). *Journal of the American Statistical Association*. 2001; 96:939–967.
- Bai J, Li K. Statistical Analysis of Factor Models of High Dimension. *The Annals of Statistics*. 2012; 40:436–465.
- Bai J, Liao Y. Statistical Inferences Using Large Estimated Covariances for Panel Data and Factor Models. *arXiv*. 2013; 1307:2662.
- Bai J, Ng S. Determining the Number of Factors in Approximate Factor Models. *Econometrica*. 2002; 70:191–221.
- Bickel PJ, Levina E. Covariance Regularization by Thresholding. *The Annals of Statistics*. 2008; 36:2577–2604.
- Cai T, Liu W. Adaptive Thresholding for Sparse Covariance Matrix Estimation. *Journal of the American Statistical Association*. 2011; 106:672–684.
- Cai T, Zhang CH, Zhou H. Optimal Rates of Convergence for Covariance Matrix Estimation. *The Annals of Statistics*. 2010; 38:2118–2144.
- Carvalho CM, Chang J, Lucas JE, Nevins JR, Wang Q, West M. High-Dimensional Sparse Factor Modeling: Applications in Gene Expression Genomics. *Journal of the American Statistical Association*. 2012; 103:1438–1456.
- Choi I. Efficient Estimation of Factor Models. *Econometric Theory*. 2012; 28:274–308.
- Fama EF, French KR. Common Risk Factors in the Returns on Stocks and Bonds. *Journal of Financial Economics*. 1993; 33:3–56.
- Fan J, Fan Y, Lv J. High Dimensional Covariance Matrix Estimation Using a Factor Model. *Journal of Econometrics*. 2008; 147:186–197.
- Fan J, Li R. Variable Selection via Nonconcave Penalized Likelihood and Its Oracle Properties. *Journal of the American Statistical Association*. 2001; 96:1348–1360.

- Fan J, Liao Y, Mincheva M. High Dimensional Covariance Matrix Estimation in Approximate Factor Models. *The Annals of Statistics*. 2011; 39:3320–3356. [PubMed: 22661790]
- Fan J, Liao Y, Mincheva M. Large Covariance Estimation by Thresholding Principal Orthogonal Complements. *Journal of the Royal Statistical, Series B*. 2013; 75:603–680.
- Johnson, RA., Wichern, DW. *Applied Multivariate Statistical Analysis*. Englewood Cliffs, NJ: Prentice Hall; 1992.
- Kapetanios G. A Testing Procedure for Determining the Number of Factors in Approximate Factor Models With Large Datasets. *Journal of Business and Economic Statistics*. 2010; 28:397–409.
- Krug T, Gabriel JP, Taipa R, Fonseca BV, Domingues-Montanari S, Fernandez-Cadenas I, Manso H, Gouveia LO, Sobral J, Albergaria I, Gaspar G, Jiménez-Conde J, Rabionet R, Ferro JM, Montaner J, Vicente AM, Rui Silva M, Matos I, Lopes G, Oliveira SA. TTC7B Emerges as a Novel Risk Factor for Ischemic Stroke Through the Convergence of Several Genome-Wide Approaches. *Journal of Cerebral Blood Flow & Metabolism*. 2012; 32:1061–1072. [PubMed: 22453632]
- Kustra R, Shioda R, Zhu M. A Factor Analysis Model for Functional Genomics. *BMC Bioinformatics*. 2006; 7:216–228. [PubMed: 16630343]
- Lam C, Fan J. Sparsistency and Rates of Convergence in Large Covariance Matrix Estimation. *The Annals of Statistics*. 2009; 37:4254–4278. [PubMed: 21132082]
- Lam C, Yao Q. Factor Modeling for High-Dimensional Time Series: Inference for the Number of Factors. *The Annals of Statistics*. 2012; 40:694–726.
- Ogata H, Goto S, Sato K, Fujibuchi W, Bono H, Kanehisa M. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids Research*. 2000; 28:27–30. [PubMed: 10592173]
- Rothman AJ, Bickel PJ, Levina E, Zhu J. Sparse Permutation Invariant Covariance Estimation. *Electronic Journal of Statistics*. 2008; 2:494–515.
- Rothman AJ, Levina E, Zhu J. Generalized Thresholding of Large Covariance Matrices. *Journal of the American Statistical Association*. 2009; 104:177–186.
- Storey JD. A Direct Approach to False Discovery Rates. *Journal of the Royal Statistical, Series B*. 2002; 64:479–498.
- Zhang CH. Nearly Unbiased Variable Selection Under Minimax Concave Penalty. *The Annals of Statistics*. 2010; 38:894–942.



**Figure 1.** Estimation error by four methods and their computational time: the dotted lines represent the means over 100 simulations and the segments represent the corresponding standard deviations.



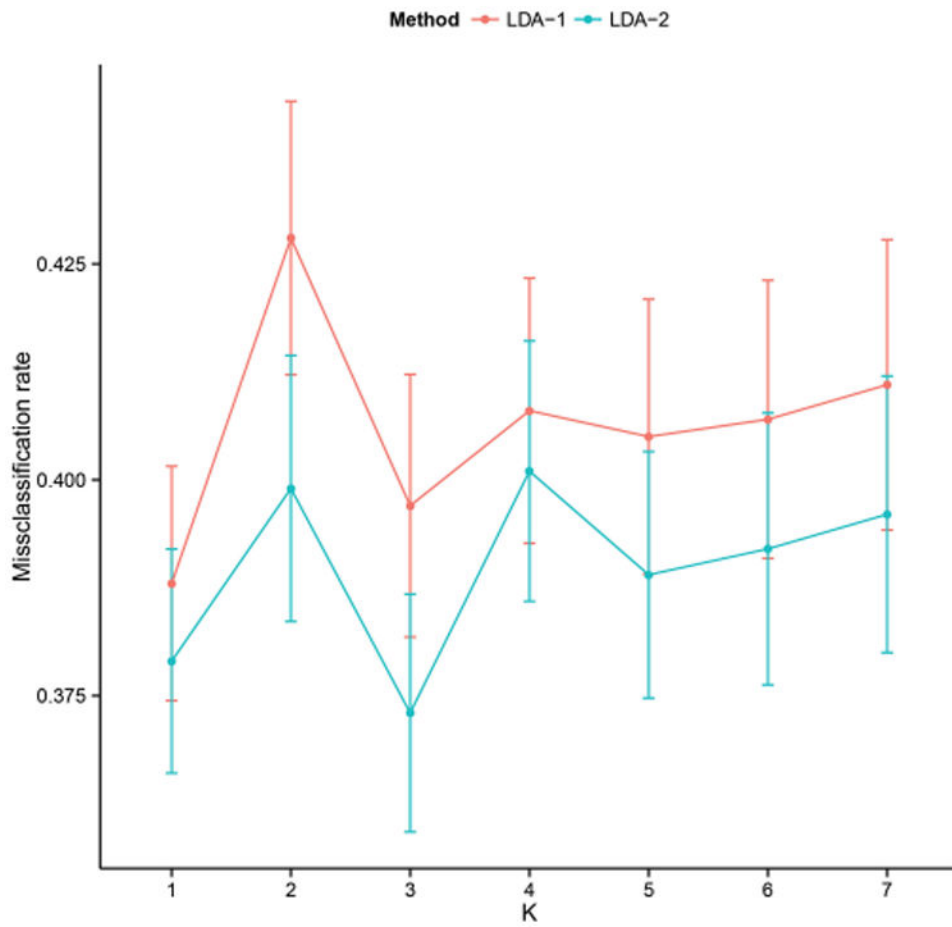
**Figure 2.** Eigen-values of the sample covariance matrix for GSE22255.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript



**Figure 3.** Misclassification rates of LDA-1 and LDA-2 over 100 random splits: the dotted lines represent the means over 100 splits and the segments represent the corresponding standard deviations.

**Table 1**

Comparison of three methods when  $s$  is much smaller than  $p(T = 200)$ .

$(s, p)$ Method	(50, 1000)			(50, 2000)		
	M1	M2	ORA	M1	M2	ORA
$\ \hat{\Sigma}_S - \Sigma_S\ _{\Sigma_S}$	0.271(0.014)	0.205(0.013)	0.204(0.013)	0.270(0.014)	0.201(0.013)	0.200(0.013)
$\ \hat{\Sigma}_S^{-1} - \Sigma_S^{-1}\ $	0.016(0.003)	0.009(0.002)	0.009(0.002)	0.017(0.003)	0.009(0.002)	0.009(0.002)
$\ \hat{\Sigma}_S - \Sigma_S\ _{\max}$	18.828(3.072)	17.460(3.237)	17.457(3.261)	18.076(2.697)	16.631(2.949)	16.623(2.950)
$\max_T \ \hat{\Gamma}_T - \mathbf{H}\ $	1.811(0.195)	0.445(0.046)	NA	1.870(0.236)	0.331(0.025)	NA
$\max_{j, s} \ \hat{\mathbf{b}}_j - \mathbf{H}\mathbf{b}_j\ $	8.064(0.694)	4.100(0.330)	3.858(0.274)	8.150(0.682)	3.932(0.292)	3.805(0.297)
$\max_{t \leq s, t \leq T} \ \hat{\mathbf{b}}_t' \hat{\mathbf{f}}_t - \mathbf{b}_t' \mathbf{f}_t\ $	11.375(1.262)	5.519(0.813)	5.268(0.843)	11.466(1.353)	5.253(0.776)	5.113(0.739)

NOTE: M1, M2, and ORA stand for Method 1, 2, and Oracle method, respectively.

**Table 2**

Comparison of three methods when  $s$  is much smaller than  $p(T = 400)$ .

$(s, p)$ Method	(50, 1000)			(50, 2000)		
	M1	M2	ORA	M1	M2	ORA
$\ \hat{\Sigma}_S - \Sigma_S\ _{\Sigma_S}$	0.186(0.009)	0.132(0.007)	0.131(0.007)	0.186(0.009)	0.131(0.008)	0.130(0.008)
$\ \hat{\Sigma}_S^{-1} - \Sigma_S^{-1}\ $	0.011(0.002)	0.004(0.001)	0.004(0.001)	0.011(0.002)	0.004(0.001)	0.004(0.001)
$\ \hat{\Sigma}_S - \Sigma_S\ _{\max}$	14.054(1.945)	11.922(2245)	11.891(2262)	14.180(2.154)	11.901(2.603)	11.900(2.604)
$\max_T \ \hat{\Gamma}_T - \mathbf{H}\Gamma\ $	1.839(0.193)	0.417(0.036)	NA	1.843(0.198)	0.305(0.026)	NA
$\max_{j, s} \ \hat{\mathbf{b}}_j - \mathbf{H}\mathbf{b}_j\ $	6.960(0.584)	2.830(0200)	2.692(0.198)	7.024(0.605)	2.761(0.188)	2.692(0.194)
$\max_{t \leq s, t \leq T} \ \hat{\mathbf{b}}_t' \hat{\mathbf{f}}_t - \mathbf{b}_t' \mathbf{f}_t\ $	11.871(1.540)	4.138(0510)	3.824(0.501)	11.457(1.569)	4.088(0.516)	3.889(0.542)

NOTE: M1, M2, and ORA stand for Method 1, 2, and Oracle method, respectively.

**Table 3**

Comparison of three methods when *sis* comparative to  $p(T = 200)$ .

$(s, p)$ Method	(800, 1000)			(800, 2000)		
	M1	M2	ORA	M1	M2	ORA
$\ \hat{\Sigma}_S - \Sigma_S\ _{\Sigma_S}$	0.440(0.006)	0.439(0.006)	0.435(0.006)	0.439(0.006)	0.436(0.006)	0.435(0.006)
$\ \hat{\Sigma}_S^{-1} - \Sigma_S^{-1}\ $	0.062(0.009)	0.062(0.009)	0.062(0.009)	0.061(0.009)	0.061(0.009)	0.062(0.012)
$\ \hat{\Sigma}_S - \Sigma_S\ _{\max}$	24565(2.626)	24562(2.609)	24567(2.599)	24511(2.883)	24543(2.847)	24536(2.851)
$\max_T \ \hat{\Gamma}_T - \mathbf{H}\Gamma\ $	0.488(0.047)	0.447(0.040)	NA	0.478(0.049)	0.337(0.038)	NA
$\max_{j, s} \ \hat{\mathbf{b}}_j - \mathbf{H}\mathbf{b}\ $	15550(0.488)	15.370(0.462)	14.418(0.271)	15595(0.551)	15.041(0.357)	14.398(0.243)
$\max_{t \leq s, t \leq T} \ \hat{\mathbf{b}}_t^* \hat{\mathbf{f}}_t - \mathbf{b}^* \mathbf{f}_t\ $	6.745(0.611)	6.680(0.635)	6.405(0.630)	6.904(0.734)	6.697(0.763)	6588(0.737)

NOTE: M1, M2, and ORA stand for Method 1, 2, and Oracle method, respectively.



**Table 4**

Comparison of three methods when  $s$  is comparative to  $p(T = 400)$ .

$(s, p)$ Method	(800, 1000)			(800, 2000)		
	M1	M2	ORA	M1	M2	ORA
$\ \hat{\Sigma}_S - \Sigma_S\ _{\Sigma_S}$	0.193(0.004)	0.192(0.004)	0.189(0.004)	0.192(0.004)	0.190(0.004)	0.188(0.004)
$\ \hat{\Sigma}_S^{-1} - \Sigma_S^{-1}\ $	0.008(0.001)	0.008(0.001)	0.008(0.001)	0.008(0.001)	0.008(0.001)	0.008(0.001)
$\ \hat{\Sigma}_S - \Sigma_S\ _{\max}$	17.062(2.603)	17.051(2.612)	17.041(2.621)	16.919(2.182)	16.891(2.206)	16.888(2.209)
$\max_T \ \hat{\Gamma}_T - \mathbf{H}\Gamma\ $	0.467(0.038)	0.423(0.036)	NA	0.466(0.038)	0.304(0.026)	NA
$\max_{j, s} \ \hat{\mathbf{b}}_j - \mathbf{H}\mathbf{b}\ $	11.009(0.298)	10.850(0.302)	10.225(0.205)	10.934(0.274)	10.530(0.213)	10.189(0.172)
$\max_{t \leq s, t \leq T} \ \hat{\mathbf{b}}_t' \hat{\mathbf{f}}_t - \mathbf{b}' \mathbf{f}_t\ $	5.367(0.577)	5.276(0.560)	4.880(0.528)	5.293(0.411)	5.024(0.461)	4.894(0.420)

NOTE: M1, M2, and ORA stand for Method 1, 2, and Oracle method, respectively.