

**ESTIMATION  
OF WEAK FACTOR MODELS**

Yoshimasa Uematsu  
Takashi Yamagata

April 2019

The Institute of Social and Economic Research  
Osaka University  
6-1 Mihogaoka, Ibaraki, Osaka 567-0047, Japan

# Estimation of Weak Factor Models

YOSHIMASA UEMATSU<sup>\*</sup> and TAKASHI YAMAGATA<sup>†</sup>

<sup>\*</sup>*Department of Economics and Management, Tohoku University*

<sup>†</sup>*Department of Economics and Related Studies, University of York*

<sup>†</sup>*Institute of Social Economic Research, Osaka University*

April 16, 2019

## Abstract

In this paper, we propose a novel consistent estimation method for the approximate factor model of Chamberlain and Rothschild (1983), with large cross-sectional and time-series dimensions ( $N$  and  $T$ , respectively). Their model assumes that the  $r$  ( $\ll N$ ) largest eigenvalues of data covariance matrix grow as  $N$  rises without specifying each diverging rate. This is weaker than the typical assumption on the recent factor models, in which all the  $r$  largest eigenvalues diverge proportionally to  $N$ , and is frequently referred to as the weak factor models. We extend the sparse orthogonal factor regression (SOFAR) proposed by Uematsu et al. (2019) to consider consistent estimation of the weak factors structure, where the  $k$ -th largest eigenvalue grows proportionally to  $N^{\alpha_k}$  with some unknown exponents  $0 < \alpha_k \leq 1$  for  $k = 1, \dots, r$ . Importantly, our method enables us to consistently estimate  $\alpha_k$  as well. In our finite sample experiment, the performance of the new estimator uniformly dominates that of the principal component (PC) estimators in terms of mean absolute loss, and its superiority gets larger as the common components become weaker. We apply our method to analyze S&P500 firm security monthly returns from January 1984 to April 2018, and the results show that the first factor is consistently near strong, whilst the second to the fourth exponents vary over months between 0.90 and 0.65 and they cross. In another application, we consider out-of-sample performance of forecasting regressions for bond yield using extracted factors by our method and by the PC, and the forecasting performance test concludes that our method outperforms the PC method.

**Keywords.** Approximate factor models, Weak factors with sparse factor loadings, Determining the number of weak factors, Non-asymptotic error bound, Factor selection consistency, Firm security returns, Forecasting bond yields.

---

<sup>\*</sup>Department of Economics and Management, Tohoku University, 27-1 Kawauchi, Aobaku, Sendai 980-8576, Japan (E-mail: yoshimasa.uematsu.e7@tohoku.ac.jp). He gratefully acknowledges the partial support of Grant-in-Aid for JSPS Overseas Research Fellow 29-60.

<sup>†</sup>Department of Economics and Related Studies, University of York, Heslington, York, YO10 5DD, UK and Institute of Social and Economic Research (ISER), Osaka University, Japan (E-mail: takashi.yamagata@york.ac.uk). He gratefully acknowledges the partial support of JSPS KAKENHI JP15H05728 and JP18K01545.

The authors appreciate Kun Chen giving helpful suggestions on computation.

## 1 Introduction

The approximate factor model with large cross-sectional and time-series dimensions ( $N$  and  $T$ , respectively) has become an increasingly important tool for the analysis of finance, macroeconomics and beyond. In finance, the model is firstly introduced by Chamberlain and Rothschild (1983), then developed in the subsequent articles by Connor and Korajczyk (1986, 1993), Bai and Ng (2002), Bai (2003), Fan et al. (2008), Fan et al. (2011, 2013), among many others. In macroeconomics, Stock and Watson (2002) propose to extract a small number of factors from the large macroeconomic and financial series and use them to forecast a macroeconomic variable of interest. Ludvigson and Ng (2009) take a similar approach to forecast bond yields. Fan et al. (2018) give an excellent review of the large dimensional factor models.

Suppose that a vector of zero-mean stationary time series  $\mathbf{x}_t \in \mathbb{R}^N$ ,  $t = 1, \dots, T$ , is generated from the factor model

$$\mathbf{x}_t = \mathbf{B}^0 \mathbf{f}_t^0 + \mathbf{e}_t, \quad (1)$$

where  $\mathbf{B}^0 = (\mathbf{b}_1^0, \dots, \mathbf{b}_N^0)' \in \mathbb{R}^{N \times r}$  with  $\mathbf{b}_i^0 \in \mathbb{R}^r$  is a matrix of deterministic factor loadings,  $\mathbf{f}_t^0 \in \mathbb{R}^r$  is a vector of probabilistic latent factors, and  $\mathbf{e}_t \in \mathbb{R}^N$  is an idiosyncratic error vector. For a while suppose  $r$  is given. Imposing the familiar identification restrictions,  $\mathbb{E}[\mathbf{f}_t^0 \mathbf{f}_t^{0'}] = \mathbf{I}_r$  and  $\mathbf{B}^0 \mathbf{B}^{0'}$  is a diagonal matrix with different elements, and assuming an exogeneity condition, we have

$$\boldsymbol{\Sigma}_x = \mathbf{B}^0 \mathbf{B}^{0'} + \boldsymbol{\Sigma}_e, \quad (2)$$

where  $\boldsymbol{\Sigma}_x = \mathbb{E}[\mathbf{x}_t \mathbf{x}_t']$  and  $\boldsymbol{\Sigma}_e = \mathbb{E}[\mathbf{e}_t \mathbf{e}_t']$ . If the eigenvalues of  $\boldsymbol{\Sigma}_e$  are uniformly bounded, the asymptotic variation of  $\boldsymbol{\Sigma}_x$  is fully controlled by loadings  $\mathbf{B}^0$ . Actually, we can observe that

$$\lambda_k(\boldsymbol{\Sigma}_x) \asymp \lambda_k(\mathbf{B}^{0'} \mathbf{B}^0)$$

for  $k = 1, \dots, r$  and  $\lambda_k(\boldsymbol{\Sigma}_x)$  are uniformly bounded for  $k = r+1, \dots, N$ . Most of the existing researches after Chamberlain and Rothschild (1983), including Connor and Korajczyk (1986, 1993), Stock and Watson (2002), Bai and Ng (2002, 2006, 2013), and Bai (2003), consider *strong factor models*, where it is assumed that all the  $r$  largest eigenvalues of  $\boldsymbol{\Sigma}_x$  diverge proportional to  $N$  (i.e.,  $\lambda_k(\mathbf{B}^{0'} \mathbf{B}^0) \asymp N$  for all  $k = 1, \dots, r$ ) while the rest of the eigenvalues remain bounded, and employ the principal components (PC) estimator.

Actually, Chamberlain and Rothschild (1983) introduced the celebrated *approximate factor models*, where  $\lambda_k(\boldsymbol{\Sigma}_x)$  diverges for each  $k = 1, \dots, r$  through assuming  $\lambda_r(\mathbf{B}^{0'} \mathbf{B}^0) \rightarrow \infty$  as  $N \rightarrow \infty$ , and the rest of the eigenvalues remain bounded (due to the uniform boundedness of  $\lambda_k(\boldsymbol{\Sigma}_e)$ ). This assumption is far less restrictive than that of the strong factor models, since the diverging rates of  $\lambda_k(\boldsymbol{\Sigma}_x)$  can be less than or equal to  $N$ , which may be different for each  $k = 1, \dots, r$ . Following De Mol et al. (2008), we call this model the *weak factor model* (WF model), in order to distinct from the strong factor (SF) model. To our knowledge, no method has been proposed which can estimate the WF models and each of the diverging rates of  $\lambda_k(\boldsymbol{\Sigma}_x)$ , jointly.

In this paper, we will fill this important gap in the literature of approximate factor models, by considering estimation of the WF models induced by *sparse* factor loadings. Specifically,  $\mathbf{B}^0$  is supposed to contain many zeros such that  $\lambda_k(\mathbf{B}^{0'} \mathbf{B}^0) \asymp N_k$ , where  $N_k = N^{\alpha_k}$  with some unknown exponent  $0 < \alpha_k \leq 1$ . Also note that our model nests the SF model when  $\alpha_k = 1$  for all  $k = 1, \dots, r$ . The sparse factor loadings are frequently observed in macroeconomic

and finance data. As an illustration, we have regressed each of 451 monthly security excess returns, which constitute the S&P500 index on December 2015, with 120 months observations back (among 500 securities) on the celebrated Fama and French (2015) five common factors, *Market*, *SMB*, *HML*, *RMW* and *CMA*, and an intercept. The numbers of securities, for which the *Market*, *SMB*, *HML*, *RMW* or *CMA* is significant at the 5% level t-test, are 446, 107, 126, 68 and 62, respectively.<sup>1</sup> Apart from the market factor, the common factors are *not* significantly different from zero for large portions of the securities. This evidence strongly suggests sparse factor loadings for the firm security returns and supports our approach<sup>2</sup>.

In practice, the number of (weak) factors,  $r$ , is unknown, so that it has to be determined in prior to the model estimation. The widely used methods for determining the number of factors proposed by Connor and Korajczyk (1993), Bai and Ng (2002), Hallin and Liška (2007), Amengual and Watson (2007), Ahn and Horenstein (2013), Caner and Han (2014), among many others, assume the SF models. The only exceptions are Onatski (2010) and Kapetanios (2010), which permit to determine the number of weak factors. One of our contributions is to provide asymptotic justification of using the edge distribution (ED) estimator by Onatski (2010) for our WF models, and to report the experimental results. The evidence therein suggests that the ED estimator is reliable for the WF models while other estimators assuming SF models are in general unreliable.

Regarding the estimation of our WF models, we no longer rely on the eigenvalue problem but on the sparsity-inducing  $\ell_1$ -norm regularization, unlike the PC estimator of the SF models. This problem requires numerical optimization, but is substantially complicated due to the imposition of both sparsity and orthogonality on the estimator. Despite this intrinsic difficulty, we propose a novel estimator of the WF models by employing the recently developed framework, the *sparse orthogonal factor regression* (SOFAR) of Uematsu et al. (2019). Hereafter the new estimator is called the WF-SOFAR estimator. As a theoretical contribution, we derive the non-asymptotic error bound for the WF-SOFAR estimator and its rate of convergence. Furthermore, we propose the *adaptive* WF-SOFAR estimator, which reduces the bias caused by the regularization and can enjoy *factor selection consistency*. This property asymptotically guarantees the true support recovery of the sparse loadings. It is remarkable that this enables us to consistently estimate each exponent  $\alpha_k$  of the divergence rates as a corollary. Since, as discussed in Bailey et al. (2016), these are interpreted as the strength of the influence of the common factors and of the cross-sectional correlations, their estimation is of great interest to empirical research.<sup>3</sup> Perhaps surprisingly, our WF-SOFAR estimator can consistently estimate the WF models with  $\alpha_k$  less than  $1/2$ . To our knowledge,

---

<sup>1</sup>Specifically, we run the time series regression  $r_{ti} - r_{ft} = a_i + b_i(r_{mt} - r_{ft}) + s_iSMB_t + h_iHML_t + r_iRMW_t + c_iCMA_t + e_{ti}$ , where  $r_{ti}$  is the  $i$ -th security monthly return at the month  $t$ ,  $r_{ft}$  is the one-month treasury bill rate,  $r_{mt}$  is the market return,  $SMB_t$  is the return on a diversified portfolio of small stocks minus the return on a diversified portfolio of big stocks,  $HML_t$  is the difference between the returns on diversified portfolios of high and low B/M stocks,  $RMW_t$  is the difference between the returns on diversified portfolios of stocks with robust and weak profitability, and  $CMA_t$  is the difference between the returns on diversified portfolios of the stocks of low and high investment firms, which is called conservative and aggressive, and  $e_{ti}$  is the error term. Then we implement the t-tests for  $b_i = 0$ ,  $s_i = 0$ ,  $h_i = 0$ ,  $r_i = 0$  and  $c_i = 0$ , referring their absolute values to 1.96. The firm security return is computed as explained in Section 6.1, and other variables are obtained from the Kenneth R. French Data Library. See Fama and French (2015) for more details of the data and the regression.

<sup>2</sup>We repeated this exercise over the window months between September 1989 and April 2018, using the 1% and 5% significance level t-tests. The full results, which suggest sparse factor loadings, are summarized and reported in the online supplement.

<sup>3</sup>Bailey et al. (2016) propose a method for estimation and inference of the exponent of the largest divergence rate,  $\alpha_1$ .

in the literature there have been no methods permitting consistent estimation of the factor model diverging at a rate slower than  $N^{1/2}$ . Note that the assumptions we will make are in line with the literature of the approximate factor models. Thus the statistical theory and the proofs we will explore are substantially different from those in Uematsu et al. (2019). In particular, the theoretical investigation of the adaptive estimator has not been considered in Uematsu et al. (2019) and is completely new to the literature.

We have to point out that the sparsity of the factor loadings is not necessarily rotation invariant. However, when the sparsity holds under the identification restrictions, each of the divergence rates of the  $r$  largest eigenvalues is invariant to the rotation of  $\mathbf{B}^0$  and  $\mathbf{f}_t^0$  with any invertible square matrix in  $\mathbb{R}^{r \times r}$ .<sup>4</sup> Finally, even when the model is not sparse at all, our WF-SOFAR estimator remains consistent to  $\mathbf{f}_t^0$  and  $\mathbf{B}^0$  because it coincide with the PC estimator, as the regularization coefficient tends to zero.

In order to evaluate the theoretical results of our WF-SOFAR estimator, we have implemented extensive finite sample experiments. In terms of the norm loss of estimators of factors, factor loadings and common components, the WF-SOFAR estimator uniformly dominates the PC estimator across all the designs we have considered. Perhaps surprisingly, the WF-SOFAR estimator of the common factors is uniformly more efficient than the PC estimator even in the most favorable experimental design of the PC, with the exponents of the divergence rates being 0.9. Also, the exponents of the divergence rates are correctly estimated, even when it is as small as 0.4, with sufficiently large sample size.

We consider two empirical applications. As the first empirical example, we apply our method to the S&P500 security excess returns, following Bailey et al. (2016). We estimate and plot the first four largest exponents of divergence rates using the 120 months estimation windows from January 1984 to April 2018. It is found that the first factor is near strong with the value of  $\alpha_1$  between 1.00 and 0.98, but the plots of the second to the fourth exponents vary over months between 0.90 and 0.65 and they also cross. A sharp rise of  $\alpha_2$  around the peak and burst of the dot-com bubble (March 2000) and a steep increase of  $\alpha_3$  at the time of the 2008 financial crisis draw our attention, which are in line with the well observed phenomenon that the correlation of the financial market tends to rise during the periods of market turmoils. As the second empirical example, we consider out-of-sample performance of forecasting regressions for bond yields using extracted factors via our method and the PC method, from a large number of macroeconomic (prediction) variables. We use analysis and data similar to that employed in Ludvigson and Ng (2009). For all the bond maturities considered, the forecasting performance test strongly rejects the null of the same forecasting performance, in favor of the alternative that our method outperforms the PC method.

We are not the first to consider the WF models. De Mol et al. (2008) consider the Bayesian forecasts for WF models and show that their rates of consistency are slower than the ones derived for the SF models based on the PC forecasts. Freyaldenhoven (2018) investigates the properties of the PC estimator in the WF models and proposes methods to estimate the number of common components diverging faster than a specific rate. Other related research includes Johnstone and Lu (2009), Onatski (2012) and Lettau and Pelger (2018), which consider the properties of the PC estimator with the bounded maximum eigenvalue of  $\Sigma_x$ , i.e.,  $\alpha_k = 0$  for all  $k$  in our WF specification.<sup>5</sup>

---

<sup>4</sup>Indeed, as can be seen in Section 6.1, in our empirical applications the sparse loadings assumption is appropriate under the identification restrictions.

<sup>5</sup>Bai and Ng (2017) focus on robust estimation of the factor models against outliers (in the sense of Candès et al. (2011)). For this purpose, they employ a ridge type estimator with  $\ell_2$  regularization, unlike our WF-SOFAR estimator which employs the  $\ell_1$  regularization. Bryzgalova (2016) considers estimation of

The rest of this paper is organized as follows. In Section 2, we formally define the WF models. Section 3 proposes the novel WF-SOFAR estimator of the WF models and considers its adaptive extension. Section 4 investigates the theoretical properties of the proposed estimators; specifically, we first consider determination of the number of weak factors and then derive the nonasymptotic error bound of the WF-SOFAR estimator and rate of convergence of the adaptive WF-SOFAR estimator. The error bound for the PC estimator in the WF models is also derived to compare with that of our estimator. Section 5 confirms the validity of our WF-SOFAR estimator in finite sample situations by Monte Carlo experiments. Section 6 gives empirical illustrations; the first example applies the WF-SOFAR to the firm security returns to uncover the market behavior and the second example reports forecasting performance of the bond yields by the WF-SOFAR. Finally, Section 7 concludes. The proofs for the main results are collected in Appendix A, and the related lemmas and their proofs as well as other supplementary materials are relegated to Appendices B-D in Online Supplements.

### 1.1 Notational remarks

For any matrix  $\mathbf{M} = (m_{ti}) \in \mathbb{R}^{T \times N}$ , we denote by  $\|\mathbf{M}\|_F$ ,  $\|\mathbf{M}\|_2$ ,  $\|\mathbf{M}\|_1$ , and  $\|\mathbf{M}\|_{\max}$  the Frobenius norm,  $\ell_2$ -induced (spectral) norm, entrywise  $\ell_1$ -norm, and entrywise  $\ell_\infty$ -norm, respectively. Specifically, they are defined by  $\|\mathbf{M}\|_F = (\sum_{t,i} m_{ti}^2)^{1/2}$ ,  $\|\mathbf{M}\|_2 = \lambda_1^{1/2}(\mathbf{M}'\mathbf{M})$ ,  $\|\mathbf{M}\|_1 = \sum_{t,i} |m_{ti}|$ , and  $\|\mathbf{M}\|_{\max} = \max_{t,i} |m_{ti}|$ , where  $\lambda_i(\mathbf{S})$  refers to the  $i$ th largest eigenvalue of any square matrix  $\mathbf{S}$ . We denote by  $\mathbf{I}_N$  and  $\mathbf{0}_{T \times N}$  the  $N \times N$  identity matrix and  $T \times N$  matrix with all the entries being zero, respectively. We use  $\lesssim$  ( $\gtrsim$ ) to represent  $\leq$  ( $\geq$ ) up to a positive constant factor. For any positive sequence  $a_n$  and  $b_n$  that converge to some points or diverge as  $n \rightarrow \infty$ , we write  $a_n \asymp b_n$  if  $a_n \lesssim b_n$  and  $a_n \gtrsim b_n$ . Moreover, we denote by  $a_n \sim b_n$  if  $a_n/b_n \rightarrow 1$ . We also use  $X \sim \mu$  to denote that random variable  $X$  has distribution  $\mu$ . For any positive values  $a$  and  $b$ ,  $a \vee b$  and  $a \wedge b$  stand for  $\max(a, b)$  and  $\min(a, b)$ , respectively. The indicator function is denoted by  $1\{\cdot\}$ .

## 2 Weak Factor Models

Consider the factor model in (1). We investigate the case of large dimensional factor models, where the number of variables  $N$  and the number of observations  $T$  diverge at the same time. For the sake of convenience, we assume the existence of some underlying divergent sequence  $n$  that satisfies the principle that  $N$  and  $T$  are both functions of  $n$  and that they simultaneously diverge as  $n \rightarrow \infty$  (i.e.,  $N = N(n) \rightarrow \infty$  and  $T = T(n) \rightarrow \infty$  as  $n \rightarrow \infty$ ). For example, we may simply suppose  $n = N \wedge T \rightarrow \infty$ . In Section 4, we also write  $T = N^\tau$  for some constant  $\tau > 0$  to understand the size of  $T$  relative to  $N$ . The number of factors  $r$  is unknown and to be determined in advance. Stacking the vectors vertically like  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_T)'$ ,  $\mathbf{F}^0 = (\mathbf{f}_1^0, \dots, \mathbf{f}_T^0)'$ , and  $\mathbf{E} = (\mathbf{e}_1, \dots, \mathbf{e}_T)'$ , we similarly rewrite model (1) as the matrix form

$$\mathbf{X} = \mathbf{F}^0 \mathbf{B}^{0'} + \mathbf{E} = \mathbf{C}^0 + \mathbf{E}, \quad (3)$$

where  $\mathbf{C}^0$  is called the matrix of common components. In order to identify  $\mathbf{F}^0$  and  $\mathbf{B}^0$  separately,  $r^2$  restrictions are required. Throughout the paper, we impose the set of identifiability conditions called PC1 by Bai and Ng (2013):  $\mathbf{F}^{0'} \mathbf{F}^0 / T = \mathbf{I}_r$  and  $\mathbf{B}^{0'} \mathbf{B}^0$  is a diagonal matrix

---

cross-sectional asset pricing models with non-diverging common components.

with distinct entries. We can see that the PC1 restriction makes the covariance matrix in (2) reduce to

$$\Sigma_x = \mathbf{B}^0 \mathbf{B}^{0'} + \Sigma_e.$$

As mentioned in the Introduction, Chamberlain and Rothschild (1983) consider approximate factor models in (3) allowing possibly different divergence rates of  $\lambda_j(\Sigma_x)$  for  $j = 1, \dots, r$  while  $\lambda_{r+1}(\Sigma_x)$  is bounded, which has recently been called the WF structure by De Mol et al. (2008). In this paper, we consider *sparsity-induced* WF models. Specifically, we assume *exactly sparse* factor loadings  $\mathbf{B}^0$  such that the sparsity of  $k$ th column (i.e., the number of nonzero elements in  $\mathbf{b}_k^0 \in \mathbb{R}^N$ ) is given by  $N_k = N^{\alpha_k}$  for  $k \in \{1, \dots, r\}$ , where  $N \geq N_1 \geq \dots \geq N_r$  (i.e.,  $1 \geq \alpha_1 \geq \dots \geq \alpha_r > 0$ ) and  $\alpha_k$ 's are unknown. Note that  $N_r$  must diverge since  $\alpha_r > 0$  and  $N \rightarrow \infty$ . We may theoretically consider the *weakly sparse* factor loadings; that is,  $\mathbf{B}^0 = (b_{ik})$  such that  $\sum_{i=1}^N |b_{ik}| \leq N_k$ . This assumption does not necessarily require exact zeros in  $\mathbf{B}^0$ . However, we choose not to pursue this direction to avoid a complicated technical issue and leave it as a question for future work.

Combining the sparsity assumption with the PC1 restriction, we then observe that there exist some constants  $d_1 > \dots > d_r > 0$  such that

$$\mathbf{B}^{0'} \mathbf{B}^0 = \text{diag}(d_1^2 N_1, \dots, d_r^2 N_r).$$

Therefore, under the assumption of uniform boundedness of  $\lambda_j(\Sigma_e)$ , it is not hard to see that

$$\lambda_j(\Sigma_x) \begin{cases} \asymp \lambda_j(\mathbf{B}^0 \mathbf{B}^{0'}) = d_j^2 N_j & \text{for } j \in \{1, \dots, r\}, \\ \text{is uniformly bounded} & \text{for } j \in \{r+1, \dots, N\}, \end{cases}$$

where the equality in the first line holds because  $\lambda_j(\mathbf{B}^0 \mathbf{B}^{0'}) = \lambda_j(\mathbf{B}^{0'} \mathbf{B}^0)$  for  $j \in \{1, \dots, r\}$ . Apparently, this specification fulfills the requirement of the WF structure.

For later use, we confirm the connection between  $\mathbf{C}^0 = \mathbf{F}^0 \mathbf{B}^{0'}$  and its singular value decomposition (SVD)  $\mathbf{C}^0 = \mathbf{U}^0 \mathbf{D}^0 \mathbf{V}^{0'}$ . Here,  $\mathbf{U}^0 \in \mathbb{R}^{T \times r}$  and  $\mathbf{V}^0 \in \mathbb{R}^{N \times r}$  are respectively matrices of the left- and sparse right-singular vectors of  $\mathbf{C}^0$  that satisfy restrictions  $\mathbf{U}^{0'} \mathbf{U}^0 / T = \mathbf{I}_r$  and  $\mathbf{V}^{0'} \mathbf{V}^0 = \mathbf{N}$  with  $\mathbf{N} = \text{diag}(N_1, \dots, N_r)$ , and  $\mathbf{D}^0 = \text{diag}(d_1, \dots, d_r) \in \mathbb{R}^{r \times r}$  is composed of the well-separated singular values  $d_1 > \dots > d_r > 0$ . The boundedness assumption on  $d_1$  makes the signal strength fully controllable through the sparsity of the loadings. In view of PC1 on model (3), it is reasonable to set  $\mathbf{F}^0 = \mathbf{U}^0$  and  $\mathbf{B}^0 = \mathbf{V}^0 \mathbf{D}^0$ . This construction yields  $\mathbf{F}^0 \mathbf{B}^{0'} = \mathbf{C}^0$  and satisfies PC1.

### 3 Estimation

In this section, we first review the PC estimator of Bai and Ng (2002) and then propose our estimator based on the SOFAR framework of Uematsu et al. (2019) for the WF models. In what follows, we refer to the WF-SOFAR estimator. In this section, we denote by  $\hat{r}$  an estimate of the number of factors. The actual method of estimating  $r$  is introduced in Section 4.1.

#### 3.1 Review of the PC estimation

In the literature of estimating model (3) with the strong factor assumption, the PC estimator has been widely used for macroeconomic and financial time series analysis. The estimator is

defined as

$$(\hat{\mathbf{F}}_{\text{PC}}, \hat{\mathbf{B}}_{\text{PC}}) = \arg \min_{(\mathbf{F}, \mathbf{B}) \in \mathbb{R}^{T \times \hat{r}} \times \mathbb{R}^{N \times \hat{r}}} \frac{1}{2} \|\mathbf{X} - \mathbf{F}\mathbf{B}'\|_{\text{F}}^2 \quad (4)$$

subject to  $\mathbf{F}'\mathbf{F}/T = \mathbf{I}_{\hat{r}}$  and  $\mathbf{B}'\mathbf{B}$  diagonal,

where  $\hat{r}$  is the predetermined number of factors. It is well-known that this optimization problem reduces to the eigenvalue problem on  $\mathbf{X}\mathbf{X}'$ ; more specifically, for given  $\hat{r}$ ,  $\hat{\mathbf{F}}_{\text{PC}}$  is obtained as  $T^{1/2}$  times the eigenvectors corresponding to the top  $\hat{r}$  largest eigenvalues of  $(NT)^{-1}\mathbf{X}\mathbf{X}'$  and  $\hat{\mathbf{B}}_{\text{PC}} = \mathbf{X}'\hat{\mathbf{F}}_{\text{PC}}/T$ . So the PC estimator does not need a special numerical optimization and is very easy to compute. Note that the PC estimator requires the SF assumption for proofs of the asymptotic normality. Under the WF assumption, however, that result is unlikely to hold.

### 3.2 WF-SOFAR estimation

Once the WF structure is well-defined via the sparsity assumption on  $\mathbf{B}^0$ , it is natural to introduce a sparsity-inducing penalty term, such as the  $\ell_1$ -norm of  $\mathbf{B}$ , to (4) to obtain a sparse estimate of  $\mathbf{B}^0$  in the same fashion as the Lasso by Tibshirani (1996). In fact, the WF-SOFAR estimator is conceptually defined as

$$(\hat{\mathbf{F}}, \hat{\mathbf{B}}) = \arg \min_{(\mathbf{F}, \mathbf{B}) \in \mathbb{R}^{T \times \hat{r}} \times \mathbb{R}^{N \times \hat{r}}} \left\{ \frac{1}{2} \|\mathbf{X} - \mathbf{F}\mathbf{B}'\|_{\text{F}}^2 + \eta_n \|\mathbf{B}\|_1 \right\} \quad (5)$$

subject to  $\mathbf{F}'\mathbf{F}/T = \mathbf{I}_{\hat{r}}$  and  $\mathbf{B}'\mathbf{B}$  diagonal.

We should note that the estimator is no longer computed by the eigenvalue problem. Even some algorithms used for the lasso estimation, such as coordinate descent, cannot be directly applied to the problem due to the seemingly incompatible restrictions, sparsity and orthogonality (diagonality).

In order to overcome this difficulty, we apply the SOFAR algorithm proposed by Uematsu et al. (2019) to solving (5). Generally speaking, the algorithm provides estimates for the SVD of a coefficient matrix in a multiple linear regression, with simultaneously exhibiting both low-rankness in the singular values matrix and sparsity in the singular vectors matrices. Recall the connection between  $(\mathbf{F}, \mathbf{B})$  and  $(\mathbf{U}, \mathbf{D}, \mathbf{V})$ , which has been defined by the SVD of  $\mathbf{C}^0$ , in Section 2. Then for given  $\hat{r}$ , the SOFAR algorithm can solve (5) to get  $(\hat{\mathbf{F}}, \hat{\mathbf{B}}) = (\hat{\mathbf{U}}, \hat{\mathbf{V}}\hat{\mathbf{D}})$ .

The algorithm to compute the WF-SOFAR estimate is based on the *augmented Lagrangian method* coupled with the *block coordinate decent*, and is numerically stable. Regularization parameter  $\eta_n$  must be determined prior to implementing the optimization. To select an “optimal”  $\eta_n$ , we can rely on cross-validation (CV) or information criteria, such as AIC, BIC, and GIC; see Uematsu et al. (2019) for more information on the computational aspects.

### 3.3 Adaptive WF-SOFAR estimation

It is interesting to observe which factors truly contribute to  $\mathbf{x}_t$  for each  $t$ . In general, the lasso estimator in a linear regression model tends to select more variables than necessary due to the bias caused by the regularization that equally penalizes the elements in the coefficient vector. To reduce the bias, Zou (2006) proposed the adaptive lasso. Here we introduce the adaptive WF-SOFAR based on a similar principle. Let  $\hat{\mathbf{B}}^{\text{ini}} = (\hat{b}_{ij}^{\text{ini}})$  denote the first-stage



initial estimator, such as the PC estimator, of  $\mathbf{B}^0$ . Then the  $(i, j)$ th element of the weighting matrix  $\mathbf{W} = (w_{ij})$  is defined as  $w_{ij} = 1/|\hat{b}_{ij}^{\text{ini}}|$ . In using this weight matrix, the adaptive WF-SOFAR estimator is defined as a minimizer of the second-stage weighted SOFAR problem:

$$(\hat{\mathbf{F}}^{\text{ada}}, \hat{\mathbf{B}}^{\text{ada}}) = \arg \min_{(\mathbf{F}, \mathbf{B}) \in \mathbb{R}^{T \times \hat{r}} \times \mathbb{R}^{N \times \hat{r}}} \left\{ \frac{1}{2} \|\mathbf{X} - \mathbf{F}\mathbf{B}'\|_{\text{F}}^2 + \eta_n \|\mathbf{W} \circ \mathbf{B}\|_1 \right\} \quad (6)$$

subject to  $\mathbf{F}'\mathbf{F}/T = \mathbf{I}_{\hat{r}}$  and  $\mathbf{B}'\mathbf{B}$  diagonal,

where  $\mathbf{A} \circ \mathbf{B}$  represents the Hadamard product of two matrices,  $\mathbf{A}$  and  $\mathbf{B}$ , of the same size. Through this construction, small weight  $w_{ij}$  indicates the importance of  $b_{ij}$ .

Recall that the  $k$ th column of  $\mathbf{B}^0$ ,  $\mathbf{b}_k^0$ , has  $N_k = N^{\alpha_k}$  nonzero entries. Similarly, let  $\hat{N}_k$  denote the number of nonzero elements in  $\hat{\mathbf{b}}_k^{\text{ada}}$ . As the Lasso in a linear regression, we may expect that the adaptive WF-SOFAR estimate  $\hat{\mathbf{B}}^{\text{ada}}$  can successfully recover the true sparsity pattern of  $\mathbf{B}^0$ . If this is true, the estimators of exponents  $\alpha_k$ 's can naturally be obtained as  $\hat{\alpha}_k = \log \hat{N}_k / \log N$  by a simple algebraic formulation. In the subsequent section, we will prove this estimator is actually consistent for  $\alpha_k$ .

## 4 Theory

In this section, we investigate the theoretical properties of the (adaptive) WF-SOFAR estimators, which are introduced in the previous section. We first reveal the asymptotic behavior of the eigenvalues of  $\mathbf{X}\mathbf{X}'$  under the weak factor structure in Section 4.1. This helps us to determine the number of weak factors. Next we derive the non-asymptotic estimation error bound for the WF-SOFAR estimator in Section 4.2. Furthermore, the asymptotic property of the adaptive WF-SOFAR estimator is derived in Section 4.3.

In order to achieve these results, we need technical assumptions on the model. Following Rigollet and Hütter (2017), we first introduce sub-Gaussian and sub-exponential random variables. A random variable  $X \in \mathbb{R}$  is said to be *sub-Gaussian* with variance proxy  $\sigma^2$  if  $\mathbb{E}[X] = 0$  and its moment generating function satisfies  $\mathbb{E}[\exp(sX)] \leq \exp(\sigma^2 s^2/2)$  for all  $s \in \mathbb{R}$ . In this case, we write  $X \sim \text{subG}(\sigma^2)$ . A random variable  $Y \in \mathbb{R}$  is said to be *sub-exponential* with parameter  $\gamma$  if  $\mathbb{E}[Y] = 0$  and its moment generating function satisfies  $\mathbb{E}[\exp(sY)] \leq \exp(\gamma^2 s^2/2)$  for all  $|s| \leq \gamma^{-1}$ . In this case, we write  $Y \sim \text{subE}(\gamma)$ . The tails of sub-Gaussian random variables decay at least as fast as the Gaussian tail. More specifically,  $X \sim \text{subG}(\sigma^2)$  has the tail probability  $\mathbb{P}(|X| > x) \leq 2 \exp\{-x^2/(2\sigma^2)\}$ , which is a crucial property in the proofs. Another important consequence is that  $\sum_{i=1}^n a_i X_i \sim \text{subG}(\sigma^2 \sum_{i=1}^n a_i^2)$  holds for  $X_i \sim \text{i.i.d. subG}(\sigma^2)$ . Note that  $\sigma^2$  is not the variance of  $X \sim \text{subG}(\sigma^2)$ . In fact, it is known that  $\mathbb{E}[X^2] \leq 4\sigma^2$ . In developing our theory, we assume sub-Gaussianity on the latent factors and idiosyncratic errors, but the condition is sometimes argued to be restrictive. In particular, a fatter-tailed distribution is preferable when financial time series are considered. In such a case, we may alternatively suppose random variable  $X$  that has the characteristic function  $\exp(-|s|^\xi)$  for some  $\xi \in (0, 2)$ . This satisfies  $\mathbb{P}(|X| > x) \lesssim c_\xi x^{-\xi}$  and, moreover,  $\mathbb{P}(\sum_{i=1}^n a_i X_i \sim \sum_{i=1}^n |a_i|^\xi c_\alpha x^{-\xi})$ . For more information, see Fan et al. (2014).

Define  $L_n = (N \vee T)^\nu - 1$  for an arbitrary positive number  $\nu$ . Throughout the paper, including all the proofs in Appendix A,  $\nu$  is assumed to be fixed. Now we are in position to introduce assumptions on the WF model.

**Assumption 1** (Latent factors). The factor matrix  $\mathbf{F}^0 = (\mathbf{f}_1^0, \dots, \mathbf{f}_T^0)'$  is specified as the vector moving average process of order  $L_n$  (VMA( $L_n$ )) such that

$$\mathbf{f}_t^0 = \sum_{\ell=0}^{L_n} \Psi_\ell \zeta_{t-\ell}, \quad \lim_{n \rightarrow \infty} \sum_{\ell=0}^{L_n} \Psi_\ell \Psi_\ell' = \mathbf{I}_r, \quad \min_{\ell \geq 0} \|\Psi_\ell\|_2 > 0,$$

where  $\zeta_t = (\zeta_{t1}, \dots, \zeta_{tr})'$  with  $\{\zeta_{tk}\}_{t,k}$  i.i.d. subG( $\sigma_\zeta^2$ ) that has  $\mathbb{E} \zeta_{tk}^2 = 1$ , and where  $\Psi_0$  is a nonsingular, lower triangular matrix.

**Assumption 2** (Factor loadings). Each column  $\mathbf{b}_k^0$  of  $\mathbf{B}^0$  has the sparsity  $N_k = N^{\alpha_k}$  such that  $0 < N_r \leq \dots \leq N_1 \leq N$  (i.e.,  $0 < \alpha_r \leq \dots \leq \alpha_1 \leq 1$ ). Furthermore,  $\mathbf{B}^{0'} \mathbf{B}^0 = \text{diag}\{d_1^2 N_1, \dots, d_r^2 N_r\}$  such that  $d_{k-1}^2 - d_k^2 \geq \delta^{1/2} d_{k-1}^2$  for  $k \in \{2, \dots, r\}$  for some positive constant  $\delta$ .

**Assumption 3** (Idiosyncratic errors). The error matrix  $\mathbf{E} = (\mathbf{e}_1, \dots, \mathbf{e}_T)'$  is specified as the VMA( $L_n$ ) such that

$$\mathbf{e}_t = \sum_{\ell=0}^{L_n} \Phi_\ell \varepsilon_{t-\ell}, \quad \limsup_{n \rightarrow \infty} \sum_{\ell=0}^{L_n} \|\Phi_\ell\|_2 < \infty,$$

where  $\varepsilon_t = (\varepsilon_{t1}, \dots, \varepsilon_{tN})'$  with  $\{\varepsilon_{ti}\}_{t,i}$  i.i.d. subG( $\sigma_\varepsilon^2$ ) and  $\Phi_0$  is a nonsingular, lower triangular matrix.

**Assumption 4** (Parameter space). The parameter space of  $\hat{\mathbf{B}}$  in optimization (5) is given by  $\mathcal{B}(\tilde{N}) = \{\mathbf{B} \in \mathbb{R}^{N \times r} : \|\mathbf{B}\|_0 \lesssim \tilde{N}/2\}$  for  $\tilde{N} \in [N_1, N]$ .

Assumptions 1 and 3 specify the stochastic processes  $\{\mathbf{f}_t\}$  and  $\{\mathbf{e}_t\}$ , respectively, to be stationary VMA( $L_n$ ), where  $L_n \sim (N \vee T)^\nu$  diverges with an arbitrary fixed positive constant  $\nu$ . This construction is regarded as the *asymptotic linear process*. Thanks to this specification, we can consider a wide range of cross-sectional and time series dependent processes for modeling  $\{\mathbf{f}_t^0\}$  and  $\{\mathbf{e}_t\}$  as asymptotic approximations of the Wold representation. By Assumption 3, we can observe that  $\lambda_1(\mathbb{E} \mathbf{e}_t \mathbf{e}_t') < \infty$ . Assumption 2 is key to our analysis and provides the sparse structure of the factor loadings  $\mathbf{B}^0$  that leads to the WF models. The sparsity makes the divergence rate of  $\lambda_k(\mathbf{B}^{0'} \mathbf{B}^0)$  possibly slower than  $N$ . This can be called *weak pervasiveness* in contrast to the so-called pervasive condition of Fan et al. (2013) that assumes the SF model (i.e.,  $N_k = N$  for all  $k \in \{1, \dots, r\}$ ). Note that PC1 is satisfied under Assumptions 1 and 2; the summability condition, together with the strong law of large numbers, gives  $\mathbf{F}^{0'} \mathbf{F}^0 / T = \mathbf{I}_r (1 + o(1))$  a.s. and the relative eigengap condition entails the eigen-separation required in  $\mathbf{B}^{0'} \mathbf{B}^0$ .

Assumption 4 is used only when the parameter estimation is considered. Note that  $\mathbf{B}^0$  is included in  $\mathcal{B}(\tilde{N})$  for any  $\tilde{N} \in [N_1, N]$  under Assumption 2. If  $\tilde{N}$  is set to  $N$ ,  $\mathcal{B}(\tilde{N})$  coincides with the whole space,  $\mathbb{R}^{T \times r}$ . Whereas, if  $\tilde{N}$  is set to  $N_1$ ,  $\mathcal{B}(\tilde{N})$  becomes as sparse as the true parameter,  $\mathbf{B}^0$ . The PC estimator always requires optimization on  $\mathcal{B}(\tilde{N})$  since it cannot be sparse, but the WF-SOFAR estimator can allow sparse  $\mathcal{B}(\tilde{N})$  with  $\tilde{N} \in [N_1, N]$  when the true loadings matrix is expected to be sparse. An important consequence of taking sparser space is that, as explained in Section 4.2, a milder condition on the WF model can be allowed. This means that the WF-SOFAR can estimate the model with a much wider class of  $(\alpha_1, \alpha_r)$  which cannot be consistently estimated by the PC. This will be one of the distinctive features of using the WF-SOFAR relative to the PC.

Under these assumptions, we obtain the following lemma which is useful for deriving the subsequent main theorems.

**Lemma 1.** *Suppose that Assumptions 1–3 hold. Then the following inequalities simultaneously hold with probability at least  $1 - O((N \vee T)^{-\nu})$ :*

- (a)  $\|\mathbf{E}\|_2 \lesssim (N \vee T)^{1/2}$ ,
- (b)  $\|\mathbf{E}\mathbf{B}^0\|_{\max} \lesssim N_1^{1/2} \log^{1/2}(N \vee T)$ ,
- (c)  $\|\mathbf{E}'\mathbf{F}^0\|_{\max} \lesssim T^{1/2} \log^{1/2}(N \vee T)$ ,
- (d)  $\max_{i \in \{1, \dots, N\}} \left| \sum_{t=1}^T (e_{ti}^2 - \mathbb{E} e_{ti}^2) \right| \lesssim T^{1/2} \log^{1/2}(N \vee T)$ .

Lemma 1 guarantees that the stochastic terms can be bounded by some deterministic sequences with high probability. As a result, in the event that the lemma occurs, we may deal with these stochastic terms as if they were deterministic sequences in the proofs. It is worth mentioning that the results are of independent interest in the literature of high-dimensional time series analysis.

#### 4.1 Determining the number of weak factors

Before investigating the properties of the estimator, we first observe the asymptotic behavior of the eigenvalues of  $\mathbf{X}\mathbf{X}'$  under the WF model. This result yields important information for determining the number of weak factors,  $r$ . Here we write  $\lambda_j$  to express  $\lambda_j((N \vee T)^{-1}\mathbf{X}\mathbf{X}')$  for notational convenience.

**Theorem 1.** *Suppose that Assumptions 1–3 hold. Then for any finite integer  $k_{\max} \geq r$  and for each  $j \in \{1, \dots, r\}$  such that*

$$\frac{N_1^{1/2} \log^{1/2}(N \vee T)}{N_j} = o(1), \quad (7)$$

*the  $j$ th largest eigenvalue of  $(N \vee T)^{-1}\mathbf{X}\mathbf{X}'$ , denoted by  $\lambda_j$ , satisfies*

$$\lambda_j \begin{cases} \geq \frac{d_j^2 N_j T}{N \vee T} (1 + o(1)) & \text{for } j \in \{1, \dots, r\}, \\ = O(1) & \text{for } j \in \{r+1, \dots, k_{\max}\}, \end{cases}$$

*with probability at least  $1 - O((N \vee T)^{-\nu})$ .*

It is worth discussing a condition that makes  $\lambda_r$  diverge for determining the number of factors. If  $\lambda_r$  is bounded, any asymptotic argument cannot distinguish it from  $\lambda_{r+1}$ , which removes the possibility of correctly detecting  $r$ . Recall that  $N_j = N^{\alpha_j}$  for some  $\alpha_j \in (0, 1]$ . Then, condition (7) for  $j = r$  is

$$\frac{N_1^{1/2} \log^{1/2}(N \vee T)}{N_r} = N^{\alpha_1/2 - \alpha_r} \log^{1/2}(N \vee T) = o(1), \quad (8)$$

which holds if  $\alpha_1/2 < \alpha_r$ . Under condition (8), Theorem 1 establishes

$$\lambda_r \gtrsim \frac{N_r T}{N \vee T} = \begin{cases} N_r T / N & \text{for } N > T, \\ N_r & \text{for } N \leq T, \end{cases}$$

with high probability. Namely,  $\lambda_r$  always diverges as long as  $N \leq T$ , but when  $N > T$  (i.e.,  $1 > \tau$  because  $T = N^\tau$ ), diverging  $\lambda_r$  requires

$$\frac{N}{N_r T} = N^{1-\alpha_r-\tau} = o(1). \quad (9)$$

This condition is rephrased as  $\alpha_r > 1 - \tau$ . When  $\lambda_r$  tends to infinity while  $\lambda_{r+1}, \dots, \lambda_{k_{\max}}$  cluster around a single point, Theorem 1 suggests the means of determining the number of weak factors,  $r$ . This presents a case in which the method of Onatski (2010) works. In that paper, the *edge distribution* (ED) estimator,

$$\hat{r}(\delta) = \max \{j = 1, \dots, k_{\max} - 1 : \lambda_j - \lambda_{j+1} \geq \delta\},$$

was proposed, where  $\delta$  is a fixed positive constant. In short, the following important corollary of Theorem 1 is obtained.

**Corollary 1.** *Suppose that Assumptions 1–3 hold. If conditions (8) and (9) are true, then  $\lambda_r$  diverges. In this case, for a fixed positive constant  $\delta$ , we have  $\hat{r}(\delta) \rightarrow r$  with probability at least  $1 - O((N \vee T)^{-\nu})$ .*

In practice,  $\delta$  should appropriately be predetermined. In fact, Onatski (2010) suggested the method based on a calibration; see that paper for full details. As long as  $\delta$  is appropriately chosen,  $\hat{r}(\delta)$  will successfully detect the true number of factors  $r$  even when the biggest gap is observed not between  $\lambda_r$  and  $\lambda_{r+1}$  but among  $\lambda_1, \dots, \lambda_r$ . Meanwhile, the method of Ahn and Horenstein (2013), which was designed for SF models, is likely to fail in detecting  $r$  in the WF models because it defines  $\hat{r}$  as the point at which the largest gap is observed among  $\lambda_1, \dots, \lambda_{k_{\max}}$ ; this is not always the case for the WF models. In Section 5, we will check the validity of Onatski’s ED estimator in our model through numerical simulations that we will compare to Ahn and Horenshtein’s *eigenvalue ratio* (ER) and *growth ratio* (GR) and Bai and Ng’s  $IC_3$  and  $BIC_3$ .

## 4.2 Non-asymptotic error bound for the WF-SOFAR estimator

Hereafter, we suppose that the WF model satisfies conditions (8) and (9) and that  $r$  is known in view of Corollary 1. We first derive the non-asymptotic error bound of our WF-SOFAR estimator.

**Theorem 2.** *Suppose that Assumptions 1–4 and conditions (8) and (9) hold. Then for*

$$\eta_n \asymp T^{1/2} \log^{1/2}(N \vee T)$$

*in optimization (5), the WF-SOFAR estimator satisfies the non-asymptotic error bound*

$$\|\hat{\mathbf{F}} - \mathbf{F}^0\|_F + \|\hat{\mathbf{B}} - \mathbf{B}^0\|_F \lesssim \frac{\kappa_n^{-1} N_1^{1/2} T^{1/2} \log^{1/2}(N \vee T)}{1 - \kappa_n^{-1} (\tilde{N} \vee T)^{1/2} \log^{1/2}(N \vee T)}$$

*with probability at least  $1 - O((N \vee T)^{-\nu})$ , where  $\kappa_n = N_r(N_r \wedge T)/N_1$ .*

Theorem 2 yields the non-asymptotic error bound that holds with high probability for any combination of  $N$  and  $T$ . It is noteworthy that Theorem 2 is not a direct consequence of the error bounds achieved by Uematsu et al. (2019). In fact, Theorem 2 does not require the assumptions of sparse factors and minimum strength of the signals that are necessary in

Uematsu et al. (2019). Moreover, we assume the stochastic factors with serial dependence and the approximate factor structure in the error term by Assumptions 1 and 3, respectively, which greatly extends the original theory of Uematsu et al. (2019).

In order to enjoy consistency,  $\kappa_n$  must diverge relatively faster than  $(\tilde{N} \vee T)^{1/2} \log^{1/2}(N \vee T)$ , where  $\tilde{N}$  is chosen from  $N_1$  to  $N$  according to the parameter space of Assumption 4. That is, we need the condition

$$\gamma_n(\tilde{N}) := \frac{N_1(\tilde{N} \vee T)^{1/2} \log^{1/2}(N \vee T)}{N_r(N_r \wedge T)} = o(1) \quad (10)$$

for assigned  $\tilde{N} \in [N_1, N]$ . With condition (10), we can obtain the rate of convergence for the WF-SOFAR estimator.

**Corollary 2.** *Suppose that Assumptions 1–4 and conditions (9) and (10) hold. Then the non-asymptotic error bound of the WF-SOFAR estimator in Theorem 2 reduces to*

$$T^{-1/2} \|\hat{\mathbf{F}} - \mathbf{F}^0\|_F = O_p \left( \frac{N_1^{1/2} \gamma_n(\tilde{N})}{(\tilde{N} \vee T)^{1/2}} \right), \quad (11)$$

$$N^{-1/2} \|\hat{\mathbf{B}} - \mathbf{B}^0\|_F = O_p \left( \frac{N_1^{1/2} T^{1/2} \gamma_n(\tilde{N})}{N^{1/2} (\tilde{N} \vee T)^{1/2}} \right). \quad (12)$$

Observe that under condition (10), rates (11) and (12) always converge to zero. It is also easily seen that the rates do not depend on  $\tilde{N}$  because  $(\tilde{N} \vee T)^{1/2}$  in the denominator and that in  $\gamma_n(\tilde{N})$  cancel out. Note that (10) with any  $\tilde{N} \in [N_1, N]$  implies (8) because  $\alpha_1 < \alpha_r + \alpha_r \wedge \tau - (\alpha_1 \vee \tau)/2 < \alpha_r + \alpha_r \wedge \tau \leq 2\alpha_r$ .

**Remark 1.** Condition (10), together with condition (9), restricts the class of WF models in response to  $\tilde{N} \in [N_1, N]$  through limiting the region of  $(\tau, \alpha_1, \alpha_r)$ , where  $\tau$  is such that  $T = N^\tau$ . To understand the meaning of condition (10), we consider the “sparsest” parameter space for  $\hat{\mathbf{B}}$ . Condition (10) with  $\tilde{N} = N_1$ , which is equivalently written as  $\alpha_1 + (\alpha_1 \vee \tau)/2 < \alpha_r + \alpha_r \wedge \tau$ , naturally brings the largest class of the WF models. Let us consider the region of  $(\tau, \alpha_1, \alpha_r)$  restricted by this condition in turn. First we find that  $\tau \in (1/2, 2]$ . The upper bound of the difference in values of  $\alpha_1$  and  $\alpha_r$  is found to be  $1/4$ , which is attainable when  $\tau \in (3/4, 1]$  and  $\alpha_1 = 1$ . Likewise, it turns out that the lower bound of  $\alpha_r$  is  $1/3$ , which is achievable when  $\alpha_1 = \alpha_r$  and  $\tau = 2/3$ . In general, the sparser the factor loadings, the smaller value of  $\tau$  is required. This is because the information on common factors can be collected only from the cross-section units with non-zero factor loadings. Hence larger cross-section units relative to  $T$  are required for the models with weaker factor structures. Contrary to the case of  $\tilde{N} = N_1$ , condition (10) with  $\tilde{N} = N$  restricts  $\alpha_r$  to be strictly larger than  $1/2$ . This is more restrictive than the case of  $\tilde{N} = N_1$  though the upper bound of the difference is the same. In sum, the WF-SOFAR can consistently estimate the WF models with exponents  $\alpha_k$ ’s smaller than or equal to  $1/2$  by supposing a sparse parameter space. The finite sample evidence in Section 5 shows that the WF-SOFAR estimator seems quite robust to the violation of the restrictions on the region of  $(\tau, \alpha_1, \alpha_r)$  discussed in this remark.

It is interesting to compare the results of the WF-SOFAR estimator obtained in Theorem 2 and Corollary 2 with those of the PC estimator for the WF models. Notice that the PC estimator must postulate the non-sparse parameter space (i.e., Assumption with  $\tilde{N} = N$ ).

**Theorem 3.** Suppose that Assumptions 1–4 with  $\tilde{N} = N$  and conditions (9) and (10) hold. Then the PC estimator satisfies the non-asymptotic error bound

$$\|\hat{\mathbf{F}}_{\text{PC}} - \mathbf{F}^0\|_{\text{F}} + \|\hat{\mathbf{B}}_{\text{PC}} - \mathbf{B}^0\|_{\text{F}} \lesssim \frac{\kappa_n^{-1} N^{1/2} T^{1/2} \log^{1/2}(N \vee T)}{1 - \kappa_n^{-1} (N \vee T)^{1/2} \log^{1/2}(N \vee T)}$$

with probability at least  $1 - O((N \vee T)^{-\nu})$ .

We can find the difference from Theorem 2 in the error bound; both  $N_1$  and  $\tilde{N}$  are replaced by  $N$ . In particular, thanks to the change in the denominator, condition (10) needs to hold with  $\tilde{N} = N$  to achieve consistency of the PC estimator. As discussed in Remark 1, this leads to a more restrictive class of WF models which can consistently be estimated.

**Corollary 3.** Suppose that Assumptions 1–4 with  $\tilde{N} = N$  and conditions (9) and (10) hold. Then the non-asymptotic error bound of the PC estimator in Theorem 3 reduces to

$$T^{-1/2} \|\hat{\mathbf{F}}_{\text{PC}} - \mathbf{F}^0\|_{\text{F}} = O_p \left( \frac{N^{1/2} \gamma_n(N)}{(N \vee T)^{1/2}} \right), \quad (13)$$

$$N^{-1/2} \|\hat{\mathbf{B}}_{\text{PC}} - \mathbf{B}^0\|_{\text{F}} = O_p \left( \frac{T^{1/2} \gamma_n(N)}{(N \vee T)^{1/2}} \right). \quad (14)$$

**Remark 2.** The potential advantage of the WF-SOFAR estimator over the PC estimator for the WF models is twofold. First, since the PC estimation does not exploit sparse parameter space (i.e.  $\tilde{N} = N$  in Assumption 4), the PC cannot consistently estimate the WF models with  $\alpha_r$  smaller than or equal to  $1/2$ , unlike the WF-SOFAR estimator (see Remark 1). Second, the rate of convergence of the PC estimator given by Corollary 3 deteriorates in comparison to that of the WF-SOFAR estimator given by Corollary 2 due to the replacement of  $N_1$  to  $N$  in the numerators. This is because the PC estimator cannot take advantage of utilizing the sparsity in the model, but the WF-SOFAR can. The finite sample evidence in Section 5 strongly supports this observation. Finally, when the model has strong factors only, namely  $N_1 = N_r = N$ , the convergence rates of the WF-SOFAR and the PC estimators become identical when  $\tilde{N} = N$ . Therefore, the WF-SOFAR estimator is likely to converge as fast as the PC estimator even when all the factors are strong.

### 4.3 Factor selection consistency of the adaptive WF-SOFAR estimator

We have derived the error bound and rate of convergence for the WF-SOFAR estimator. We next investigate the asymptotic property of the adaptive WF-SOFAR estimator introduced in Section 3.3. Specifically, we prove the *factor selection consistency*, which guarantees that the adaptive WF-SOFAR can recover the true sparsity pattern of the loadings and correctly select the relevant factors. As a corollary of the factor selection consistency, we finally establish the consistency of the estimated exponents of the divergence rates.

Before stating the theorem, define the index set of nonzero signals in  $\mathbf{B}^0$  as

$$\mathcal{S} = \text{supp}(\mathbf{B}^0) \subset \{1, \dots, N\} \times \{1, \dots, r\}.$$

For any (sparse) matrix  $\mathbf{A} = (a_{ik}) \in \mathbb{R}^{N \times r}$ , we define  $\mathbf{A}_{\mathcal{S}} = (a_{ik} 1_{\{(i,k) \in \mathcal{S}\}})$  and  $\mathbf{a}_{\mathcal{S}} = \text{vec } \mathbf{A}_{\mathcal{S}} \in \mathbb{R}^{rN}$ . We further write  $\underline{b}_n^0 = \min_{(i,k) \in \mathcal{S}} |b_{ik}^0|$  and introduce condition

$$\frac{N_1^2}{N_r^2 T^{1/2}} = N^{2\alpha_1 - 2\alpha_r - \tau/2} = o(1). \quad (15)$$

Condition (15) further restricts the region of  $(\tau, \alpha_1, \alpha_r)$  in terms of the maximum difference of  $\alpha_1$  and  $\alpha_r$  when  $\tau < 1$ . However, the difference can be  $1/4$ , which is the same as the case constrained only by (10), as long as  $\tau = 1$ . The lower bound of  $\alpha_r$  can also achieve  $1/3$  even if (15) is additionally supposed.

**Theorem 4.** *Suppose that Assumptions 1–3 and conditions (9), (10) and (15) hold and that the weighting matrix  $\mathbf{W}$  is constructed by an estimator  $\hat{\mathbf{B}}^{\text{ini}}$  such that*

$$\|\hat{\mathbf{B}}^{\text{ini}} - \mathbf{B}^0\|_{\max} \lesssim \frac{N^{1/2}\gamma_n(\tilde{N})}{(\tilde{N} \vee T)^{1/2}} \quad (16)$$

with probability at least  $1 - O((N \vee T)^{-\nu})$ . If  $\underline{b}_n^0$  and  $\eta_n$  satisfy the conditions

$$\underline{b}_n^0 \gtrsim \frac{N_r(N_r \wedge T)\eta_n}{N_1(N_r \vee T)T^{1/2} \log^{1/2}(N \vee T)} \vee \frac{N^{1/2}\gamma_n(\tilde{N})}{(\tilde{N} \vee T)^{1/2}}, \quad (17)$$

$$\eta_n \asymp \frac{N^{1/2}T^{1/2}}{(\tilde{N} \vee T)^{1/2}}, \quad (18)$$

then the adaptive WF-SOFAR estimator satisfies

$$T^{-1/2} \left\| \hat{\mathbf{F}}^{\text{ada}} - \mathbf{F}^0 \right\|_{\text{F}} = O_p \left( \frac{N_1^{1/2}\gamma_n(\tilde{N})}{(\tilde{N} \vee T)^{1/2}} \right), \quad (19)$$

$$N^{-1/2} \left\| \hat{\mathbf{B}}_{\mathcal{S}}^{\text{ada}} - \mathbf{B}_{\mathcal{S}}^0 \right\|_{\text{F}} = O_p \left( \frac{N_1^{1/2} T^{1/2}\gamma_n(\tilde{N})}{N^{1/2}(\tilde{N} \vee T)^{1/2}} \right), \quad (20)$$

$$\mathbb{P} \left( \text{supp}(\hat{\mathbf{B}}^{\text{ada}}) = \mathcal{S} \right) \rightarrow 1. \quad (21)$$

As described in Section 3.1, we can use the PC estimator as the initial estimator. Lemma 6 in Appendix A reveals that the PC estimator satisfies condition (16). Condition (17) restricts the magnitude of the minimum signal  $\underline{b}_n^0 = \min_{(i,k) \in \mathcal{S}} |b_{ik}^0|$  from below. In the literature of sparse estimation, this kind of condition is said to be the *beta-min* condition and is frequently supposed to achieve variable selection consistency (e.g., Fan and Lv, 2011). Notice that the lower bound always converges to zero under condition (10) even in the worst case. As a result, condition (17) can asymptotically allow a wide class of the WF models. The rates of convergence (19) and (20) in Theorem 4 are identical to (11) and (12) in Corollary 2, respectively, hence they converge to zero.

Next, we prove that  $\hat{\alpha}_k = \log \hat{N}_k / \log N$ , which is defined in Section 3.3, is consistent to  $\alpha_k$ . Due to the factor selection consistency of (21) in Theorem 4, we immediately reach to the following corollary.

**Corollary 4.** *If the model selection consistency in (21) holds, then we have*

$$\mathbb{P}(\hat{\alpha}_k = \alpha_k \text{ for all } k = 1, \dots, r) \rightarrow 1.$$

In the literature on the adaptive Lasso and penalized regressions with folded-concave penalties, such as the SCAD by Fan and Li (2001) or the MCP by Zhang (2014), the asymptotic normality can be established for the nonzero subvector of the estimator as well as the variable selection consistency; this is known as the *oracle property*. It was thought to be useful for statistical inference for the estimated model, but has been criticized recently by

Leeb and Pötscher (2005, 2006, 2008) and Pötscher and Leeb (2009), among others, on the basis that inferences based on the property lack uniformity over sequences of models that include even minor deviations from the beta-min condition. The same criticism could apply to the WF models and the adaptive WF-SOFAR estimator, and hence we do not consider inference based on the adaptive WF-SOFAR estimator in the present paper.

Instead of the adaptive estimation and the *oracle property*, a number of methods have been proposed for inference in high-dimensional linear regressions. Especially, the method called *debiasing (desparsification)* by Javanmard and Montanari (2014), van de Geer et al. (2014), and Zhang and Zhang (2014) has gained popularity. This framework tries directly to correct the Lasso estimator to remove the bias using the Karush-Kuhn-Tucker (KKT) conditions which the Lasso must satisfy. Such debiased Lasso estimators for time series models have been proposed by, for instance, Kock and Tang (2019). However, the question of how to desparsify the WF-SOFAR estimator for the WF models is nontrivial. We leave the problem as a future challenge.

## 5 Monte Carlo Experiments

In this section we investigate the finite sample behavior of estimators of the number of factors, and the behavior of the proposed WF-SOFAR estimators by means of Monte Carlo experiments.

In this section, indexes  $i$ ,  $t$ , and  $k$  run over  $1, \dots, N$ ,  $1, \dots, T$ , and  $1, \dots, r$ , respectively, unless otherwise noted. Denote by  $N_k = \lfloor N^{\alpha_k} \rfloor$ , where  $\lfloor \cdot \rfloor$  is the floor function, with  $0 < \alpha_k \leq 1$  for each  $k$ . We consider the following Data Generating Process (DGP):

$$x_{ti} = \sum_{k=1}^r b_{ik} f_{tk} + \sqrt{\theta} e_{ti}. \quad (22)$$

The factor loadings  $b_{ik}$  and factors  $f_{tk}$  are formed such that  $N^{-1} \sum_{i=1}^N b_{ik} b_{i\ell} = 1\{k = \ell\}$  and  $T^{-1} \sum_{t=1}^T f_{tk} f_{t\ell} = 1\{k = \ell\}$ , by applying Gram-Schmidt orthonormalization to  $b_{ik}^*$  and  $f_{tk}^*$ , respectively, where  $b_{ik}^* \sim \text{i.i.d.} N(0, 1)$  for  $i = 1, \dots, N_k$  and  $b_{ik}^* = 0$  for  $i = N_k + 1, \dots, N$ , and

$$f_{tk}^* = \rho_{fk} f_{t-1,k}^* + v_{tk}$$

with  $v_{kt} \sim \text{i.i.d.} N(0, 1 - \rho_{fk}^2)$  and  $f_{0k}^* \sim \text{i.i.d.} N(0, 1)$ . The idiosyncratic errors  $e_{ti}$  are generated by

$$e_{ti} = \rho_e e_{t-1,i} + \beta \varepsilon_{t,i-1} + \beta \varepsilon_{t,i+1} + \varepsilon_{ti},$$

where  $\varepsilon_{ti} \sim \text{i.i.d.} N(0, \sigma_{\varepsilon,ti}^2)$  with  $\sigma_{\varepsilon,ti}^2$  being set such that  $\text{Var}(e_{ti}) = 1$ .

The above DGP is in line with that considered in the existing representative literature of approximate factor models, such as Bai and Ng (2002), Onatski (2010), and Ahn and Horenstein (2013), among many others. The main difference in our DGP from the literature is that the absolute sums of the factor loadings over  $i$  are allowed to diverge proportionally to  $N_k$ .<sup>6</sup>

As the benchmark DGP, we set  $r = 2$ ,  $\rho_{fk} = \rho_e = 0.5$  for all  $k$ ,  $\beta = 0.2$ , and  $\theta = 1$ . We focus on the performance of the estimators for different values of exponents  $(\alpha_1, \alpha_2)$ . In

---

<sup>6</sup>Ahn and Horenstein (2013) allow the number of pairs of cross-correlated errors,  $e_{ti}$ , to rise with  $N$ . The approximate factor models (of Chamberlain and Rothschild (1983)) are not in line with such a design.



particular, we consider the combinations (0.9, 0.9), (0.9, 0.8), (0.9, 0.5), (0.8, 0.8), (0.8, 0.5), and more challenging cases, (0.8, 0.4), (0.5, 0.5) and (0.5, 0.4). All the experimental results are based on 1,000 replications.

### 5.1 Determining the number of weak factors

As in Section 4.1, we write  $\lambda_k$  to express  $\lambda_k((N \vee T)^{-1}\mathbf{X}\mathbf{X}')$ . As discussed in that section,  $\lambda_r$  will diverge while  $\lambda_{r+1}$  will be bounded. Due to this, as summarized in Corollary 1, we can consider a class of estimators

$$\hat{r}(\delta) = \max \{j = 1, \dots, k_{\max} - 1 : \lambda_j - \lambda_{j+1} \geq \delta\},$$

where  $k_{\max} (\geq r)$  is given. The estimator is the maximum value of  $k$  with which  $\lambda_k - \lambda_{k+1}$  exceeds the threshold  $\delta$ . The question is how to choose the threshold value  $\delta$ . We employ ED (edge distribution) algorithm to calibrate  $\delta$ , which is proposed by Onatski (2010; p.1008). The algorithm exploits the square root shape of the edge of the eigenvalue distribution.<sup>7</sup> We denote the associated estimator of  $r$  as

$$\hat{r}_{\text{ED}} = \hat{r}(\hat{\delta}) \quad (23)$$

where  $\hat{\delta}$  is the calibrated  $\delta$ .

We consider other competitor statistics to determine the number of factors. The ER (eigenvalue ratio) and GR (growth ratio) estimators of Ahn and Horenstein (2013) are closely related to Onatski (2010) in that they utilize the eigenvalues. Under the classical approximate factor model, the probability limit of  $N^{-1} \sum_{i=1}^N (b_{ki}^0)^2$  is a bounded positive constant for all  $k$ , while the cross correlation among the idiosyncratic errors is typically bounded in  $N$ . Therefore, it is sufficient to distinguish a single big “jump” in the value of several largest eigenvalues. We define the ER and GR estimators as the maximizers of the following functions of the eigenvalues:

$$\hat{r}_{\text{ER}} = \max_{1 \leq k \leq k_{\max}} ER(k), \quad ER(k) = \frac{\lambda_k}{\lambda_{k+1}}, \quad (24)$$

$$\hat{r}_{\text{GR}} = \max_{1 \leq k \leq k_{\max}} GR(k), \quad GR(k) = \frac{\log[V(k-1)/V(k)]}{\log[V(k)/V(k+1)]} \quad (25)$$

with  $V(k) = \sum_{\ell=k+1}^{N \wedge T - k_{\max} - 1} \lambda_{\ell}$ . Note that, as these estimators just identify the maximum gap between the ordered eigenvalues, when the gap of divergence rate of factors,  $\alpha_k - \alpha_{k+1}$ , is relatively large, these statistics might pick up  $k$  as the estimator of  $r$ , even when  $k < r$ .

The information criteria proposed by Bai and Ng (2002) are widely used to determine the number of factors. Among these criteria,  $IC_3$  and  $BIC_3$  are of relevance for the models in which we are interested.  $IC_3$  is appropriate for large  $N$  and  $T$  panel data and also recommended by Bai and Ng (2002).  $BIC_3$  is recommended by Bai to include in the simulation exercises of Ahn and Horenstein (2013), particularly for cases in which idiosyncratic errors are cross-sectionally correlated. We define the associated estimators

$$\hat{r}_{\text{IC}_3} = \min_{1 \leq k \leq k_{\max}} IC_3(k), \quad \hat{r}_{\text{BIC}_3} = \min_{1 \leq k \leq k_{\max}} BIC_3(k),$$

---

<sup>7</sup>We have found no experimental results on the finite sample performance of the ED estimator with the WF models apart from ours.

where

$$\begin{aligned} IC_3(k) &= \log v(k) + \frac{k \log(N \wedge T)}{N \wedge T}, \\ BIC_3(k) &= v(k) + \frac{kv(k_{\max})(N + T - k) \log(NT)}{NT}, \end{aligned} \quad (26)$$

with

$$v(k) = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \left( x_{ti} - \sum_{\ell=1}^k \hat{b}_{PCi\ell} \hat{f}_{PCt\ell} \right).$$

Observe that computation of the information criteria requires estimation of the factor models, which is not necessary for  $\hat{r}_{ED}$ ,  $\hat{r}_{ER}$ , and  $\hat{r}_{GR}$ .

### 5.1.1 Results

Table 1 reports the average of the estimated number of factors over the replications by the ED of Onatski (2010), GR of Ahn and Horenstein (2013), and  $BIC_3$  of Bai and Ng (2002), which are defined by (23), (25) and (26), respectively.<sup>8</sup> We employ the benchmark DGP of equation (22) with  $r = 2$ ,  $\rho_{fk} = \rho_e = 0.5$  for  $k = 1, \dots, r$ ,  $\beta = 0.2$ ,  $\theta = 1$  and the sets of the values of exponents  $(\alpha_1, \alpha_2)$ , as described above. We set the maximum number of factors,  $k_{\max}$ , as five. All the combinations of  $N, T = 100, 200, 500, 1000$  are considered and all the results are based on 1,000 replications. As can be seen in Table 1, when  $\alpha_1$  and  $\alpha_2$  are both close to unity, all the methods perform very well, picking up the true number of factors with very high probability. Indeed, in the case of exponents  $(\alpha_1, \alpha_2) = (0.9, 0.9), (0.9, 0.8), (0.8, 0.8)$ , GR and  $BIC_3$  choose the correct number of factors for all the replications, whilst ED very slightly tends to overestimate the number of factors.

However, the performance of GR and  $BIC_3$  deteriorates when the gap of the values between  $\alpha_1$  and  $\alpha_2$  widens, or when both values  $\alpha_1$  and  $\alpha_2$  are further away from unity. For example, when  $(\alpha_1, \alpha_2) = (0.9, 0.5)$ , for  $N = T = 100, 200, 500, 1000$ , the averages of the estimated numbers of factors by GR are 1.15, 1.10, 1.05 and 1.00, respectively. It converges to one as the sample size increases, despite the fact that the true number of factors is two. With the same sets of exponents, for  $N = T = 100, 200, 500, 1000$ , the averages of the estimated number of factors by  $BIC_3$  are 1.30, 1.36, 1.53 and 1.42, which apparently does not tend towards the true number of factors as the sample size increases. In contrast, ED performs very well, and its estimation quality is very similar to that when both exponents are close to unity. Similar observations for ED, GR and  $BIC_3$  apply to the cases in which  $(\alpha_1, \alpha_2) = (0.8, 0.5), (0.8, 0.4)$ , however, for the latter case ED tends to underestimate the number of factors when  $T < 500$  and  $N < 200$ . When  $(\alpha_1, \alpha_2) = (0.5, 0.5), (0.5, 0.4)$ , as the values of  $\alpha_1$  and  $\alpha_2$  are similar, GR might be expected to pick up the right number of factors, but it tends to significantly underestimate  $r$ . For example, when  $(\alpha_1, \alpha_2) = (0.5, 0.4)$ , for  $N = T = 1000$ , the average of the estimated number of factors is 1.59.  $BIC_3$  performs even worse, giving the values between 1.00 and 1.23 for all the combinations of  $N$  and  $T$ . Under this challenging set up, ED consistently estimates the number of factors for sufficiently large  $T$  and  $N$ . For example, when  $(\alpha_1, \alpha_2) = (0.5, 0.4)$ , for  $N = T = 500, 1000$ , the averages of the estimated numbers of factors are 2.01 and 2.02.

<sup>8</sup>To save the space, we do not report the results for ER and  $IC_3$  since the performance of ER is very similar to that of GR, and the performance of  $IC_3$  is mostly outperformed by  $BIC_3$ . These results are available upon request from the authors.

We conclude that the finite sample evidence suggests that the ED method of Onatski (2010) provides a reliable estimation of the number of factors in approximate factor models, whilst the methods proposed by Ahn and Horenstein (2013) and Bai and Ng (2002) may not be as reliable as the ED, in general.

Table 1: Average of the chosen number of factors for WF models by edge distribution algorithm (ED), Growth Ratio (GR), and BIC<sub>3</sub> methods:  $r = 2$ ,  $k_{\max} = 5$

$T, N$	ED				GR				BIC3			
	100	200	500	1000	100	200	500	1000	100	200	500	1000
$(\alpha_1, \alpha_2) = (0.9, 0.9)$												
100	2.05	2.04	2.02	2.01	2.00	2.00	2.00	2.00	2.00	2.00	2.00	2.00
200	2.04	2.04	2.03	2.02	2.00	2.00	2.00	2.00	2.00	2.00	2.00	2.00
500	2.04	2.04	2.03	2.02	2.00	2.00	2.00	2.00	2.00	2.00	2.00	2.00
1000	2.02	2.04	2.03	2.02	2.00	2.00	2.00	2.00	2.00	2.00	2.00	2.00
$(\alpha_1, \alpha_2) = (0.9, 0.8)$												
100	2.04	2.03	2.02	2.01	2.00	2.00	2.00	2.00	2.00	2.00	2.00	2.00
200	2.04	2.03	2.03	2.02	2.00	2.00	2.00	2.00	2.00	2.00	2.00	2.00
500	2.03	2.03	2.02	2.02	2.00	2.00	2.00	2.00	2.00	2.00	2.00	2.00
1000	2.04	2.03	2.02	2.02	2.00	2.00	2.00	2.00	2.00	2.00	2.00	2.00
$(\alpha_1, \alpha_2) = (0.9, 0.5)$												
100	1.96	1.95	1.95	1.90	1.15	1.05	1.00	1.00	1.30	1.17	1.02	1.00
200	2.02	2.02	2.03	2.02	1.23	1.10	1.01	1.00	1.40	1.36	1.12	1.01
500	2.04	2.03	2.02	2.02	1.41	1.22	1.05	1.00	1.41	1.51	1.53	1.42
1000	2.02	2.03	2.02	2.02	1.35	1.22	1.05	1.00	1.43	1.51	1.53	1.42
$(\alpha_1, \alpha_2) = (0.8, 0.8)$												
100	2.05	2.03	2.02	2.01	2.00	2.00	2.00	2.00	2.00	2.00	2.00	2.00
200	2.05	2.03	2.03	2.02	2.00	2.00	2.00	2.00	2.00	2.00	2.00	2.00
500	2.03	2.03	2.02	2.02	2.00	2.00	2.00	2.00	2.00	2.00	2.00	2.00
1000	2.03	2.03	2.02	2.02	2.00	2.00	2.00	2.00	2.00	2.00	2.00	2.00
$(\alpha_1, \alpha_2) = (0.8, 0.5)$												
100	1.96	1.96	1.95	1.90	1.30	1.18	1.04	1.00	1.30	1.17	1.02	1.00
200	2.02	2.02	2.03	2.02	1.40	1.30	1.09	1.01	1.39	1.36	1.12	1.01
500	2.03	2.03	2.02	2.02	1.61	1.45	1.24	1.10	1.41	1.51	1.53	1.42
1000	2.02	2.03	2.02	2.02	1.52	1.45	1.24	1.10	1.43	1.51	1.53	1.42
$(\alpha_1, \alpha_2) = (0.8, 0.4)$												
100	1.70	1.66	1.52	1.35	1.09	1.03	1.00	1.00	1.09	1.02	1.00	1.00
200	1.86	1.90	1.89	1.86	1.12	1.05	1.00	1.00	1.10	1.06	1.00	1.00
500	2.01	2.00	2.01	2.02	1.21	1.11	1.01	1.00	1.10	1.13	1.05	1.01
1000	1.94	2.00	2.01	2.02	1.18	1.11	1.01	1.00	1.12	1.13	1.05	1.01
$(\alpha_1, \alpha_2) = (0.5, 0.5)$												
100	1.80	1.84	1.79	1.72	1.64	1.62	1.65	1.62	1.08	1.01	1.00	1.00
200	1.98	2.02	2.02	2.02	1.69	1.80	1.83	1.87	1.11	1.09	1.01	1.00
500	2.03	2.03	2.02	2.02	1.87	1.91	1.95	1.98	1.14	1.21	1.23	1.14
1000	2.02	2.03	2.02	2.02	1.78	1.91	1.95	1.98	1.14	1.21	1.23	1.14
$(\alpha_1, \alpha_2) = (0.5, 0.4)$												
100	1.54	1.52	1.36	1.14	1.50	1.47	1.39	1.33	1.03	1.00	1.00	1.00
200	1.83	1.88	1.89	1.86	1.52	1.53	1.50	1.39	1.03	1.02	1.00	1.00
500	2.00	2.00	2.01	2.02	1.67	1.64	1.65	1.59	1.03	1.05	1.02	1.01
1000	1.92	2.00	2.01	2.02	1.60	1.64	1.65	1.59	1.04	1.05	1.02	1.01

Notes: The DGP is  $x_{ti} = \sum_{k=1}^r b_{ik}^0 f_{tk}^0 + \sqrt{\theta} e_{ti}$ ,  $i = 1, \dots, N$ ,  $t = 1, \dots, T$ , where  $b_{ik}^* \sim \text{i.i.d.} N(0, 1)$  for  $i = 1, \dots, \lfloor N^{\alpha_k} \rfloor = N_k$  and  $b_{ik}^* = 0$  for  $i = N_k + 1, \dots, N$ , with  $0 < \alpha_k \leq 1$ ,  $f_{tk}^* = \rho_{fk} f_{k,t-1}^* + v_{tk}$  for  $t = 1, \dots, T$ ,  $v_{tk} \sim \text{i.i.d.} N(0, 1 - \rho_{fk}^2)$  with  $f_{0k}^* \sim \text{i.i.d.} N(0, 1)$  for  $k = 1, \dots, r$ ,  $e_{ti} = \rho_e e_{t-1,i} + \beta \varepsilon_{t,i-1} + \beta \varepsilon_{t,i+1} + \varepsilon_{ti}$ ,

$\varepsilon_{ti} \sim \text{i.i.d.} N(0, \sigma_{\varepsilon, ti}^2)$ ,  $\sigma_{\varepsilon, ti}^2$  is set so that  $\text{Var}(e_{ti}) = 1$ . Once  $b_{ik}^*$  and  $f_{tk}^*$  are generated, we apply Gram-Schmidt orthonormalization and form  $b_{ik}^0$  and  $f_{tk}^0$ , where  $N^{-1} \sum_{i=1}^N b_{ik}^0 b_{i\ell}^0 = 1\{k = \ell\}$  and  $T^{-1} \sum_{t=1}^T f_{tk}^0 f_{t\ell}^0 = 1\{k = \ell\}$ . We set  $r = 2$ ,  $\rho_{fk} = \rho_e = 0.5$  for  $k = 1, \dots, r$ ,  $\beta = 0.2$ ,  $\theta = 1$ . The reported values are the average of estimated number of factors over the replications by the edge distribution algorithm (ED) of Onatski (2010), growth ratio (GR) proposed by Ahn and Horenstein (2013), and  $BIC_3$  of Bai and Ng (2002), which are defined by (23), (25) and (26), respectively. All the results are based on 1,000 replications.

## 5.2 Finite sample properties of the WF-SOFAR estimator

In this subsection we investigate the finite sample properties of our WF-SOFAR estimator, and compare these with those of the PC estimator where appropriate. Here we treat the number of factors,  $r$ , as given, in order to distinguish the analysis from the quality of methods for choosing the number of factors, which is investigated above. We employ the benchmark DGP of equation (22) with  $r = 2$ ,  $\rho_{fk} = \rho_e = 0.5$  for  $k = 1, \dots, r$ ,  $\beta = 0.2$ ,  $\theta = 1$ . Initially we consider the exponents  $(\alpha_1, \alpha_2) = (0.9, 0.9)$ ,  $(0.9, 0.8)$ ,  $(0.9, 0.5)$ ,  $(0.8, 0.8)$ ,  $(0.8, 0.5)$  and all the combinations of  $N, T = 100, 200, 500, 1000$ .<sup>9</sup> Then, we investigate more challenging cases,  $(0.8, 0.4)$ ,  $(0.5, 0.5)$  and  $(0.5, 0.4)$ . In these cases, we consider large sample sizes,  $N, T = 500, 1000$ , only. All the results are based on 1000 replications.

We report the results of the adaptive WF-SOFAR estimator with BIC, which we recommend to use.<sup>10</sup> We consider the PC estimator as defined in Section 3.1.

For performance comparison purposes, we consider the following  $\ell_2$ -norm losses based on the  $k$ th scaled estimators:

$$\begin{aligned} L(\widehat{\mathbf{F}}) &= \left\| \sum_{k=1}^r T^{-1/2} [\text{abs}(\widehat{\mathbf{f}}_k) - \text{abs}(\mathbf{f}_k^0)] \right\|_2, \\ L(\widehat{\mathbf{B}}) &= \left\| \sum_{k=1}^r N_k^{-1/2} [\text{abs}(\widehat{\mathbf{b}}_k) - \text{abs}(\mathbf{b}_k^0)] \right\|_2, \\ L(\widehat{\mathbf{C}}) &= \left\| \sum_{k=1}^r T^{-1/2} N_k^{-1/2} [\widehat{\mathbf{C}}_k - \mathbf{C}_k^0] \right\|_F, \end{aligned} \tag{27}$$

where  $\text{abs}(\mathbf{a})$  takes elementwise absolute value of a real vector  $\mathbf{a}$ . Due to the scaling, the performance of the estimators can be comparable across different combinations of the values of  $N$ ,  $T$ , and  $\alpha_k$ 's. Observe that these norm losses are not sensitive to the sign indeterminacy of the estimators (i.e.  $\mathbf{f}_k^0 \mathbf{b}_k^{0'} = (-\mathbf{f}_k^0)(-\mathbf{b}_k^{0'})$ ) and the change of the order of the factor components (e.g., for  $r = 2$ , the estimated first factor can be of the true second factor).

### 5.2.1 Results

Table 2 reports the averages and standard deviations (s.d.) of the estimates of  $\alpha_1$  and  $\alpha_2$  based on Corollary 4, and the average of the norm losses (multiplied by 100) of the scaled estimated factors, factor loadings, and common components by the WF-SOFAR (WS in the tables) and PC estimators over the replications. The experimental design is essentially

<sup>9</sup>Note that when the values of  $\alpha_1$  are 0.9 and 0.8, the associated lowest bounds of  $\alpha_r$  implied by condition (10) are 0.675 and 0.6, respectively.

<sup>10</sup>We examined all the combinations of WF-SOFAR and adaptive WF-SOFAR with AIC, cross-validation, BIC and GIC. The results of which are available upon request from the authors.

determined by the two values of the exponents,  $(\alpha_1, \alpha_2)$ . In this table, we consider  $(\alpha_1, \alpha_2) = (0.9, 0.9), (0.9, 0.8), (0.9, 0.5), (0.8, 0.8), (0.8, 0.5)$ . Each panel of the table has two column blocks for  $T = 100$  and  $200$  (and  $T = 500$  and  $1000$ ), and each column block for given  $T$  contains four row blocks for  $N = 100, 200, 500, 1000$ .

First, let us focus on the WF-SOFAR estimates of  $(\alpha_1, \alpha_2)$ . In a nutshell, they are sufficiently accurate but tend to slightly underestimate when  $\alpha_k$  is closer to one and overestimate when it is around 0.5. The precision improves as  $T$  and  $N$  increase. For example, when  $(\alpha_1, \alpha_2) = (0.9, 0.5)$ , for  $N = T = 100, 200, 500, 1000$ , the averages of  $(\hat{\alpha}_1, \hat{\alpha}_2)$  are  $(0.85, 0.54), (0.87, 0.53), (0.88, 0.52), (0.89, 0.52)$  and the standard deviations (rounded to two decimal places) are  $(0.02, 0.07), (0.01, 0.04), (0.00, 0.03), (0.00, 0.02)$ . This precision is remarkable since given  $\alpha_1$  is 0.9, the value of the lower bound of  $\alpha_r$  implied by condition (10) is 0.675, which is much larger than the actual value considered here, 0.5. Similar comments apply to the results for the combinations  $(\alpha_1, \alpha_2) = (0.8, 0.5)$ .

Now we turn our attention to the performance of the WF-SOFAR and PC estimates. In terms of the norm loss in (27), the WF-SOFAR uniformly beats the PC across all the designs we have considered. Perhaps surprisingly, the WF-SOFAR estimate of the factors is much more accurate than the PC estimator even in the most favorable experimental design to the PC, with  $(\alpha_1, \alpha_2) = (0.9, 0.9)$ . For example, for  $N = 100, 200, 500, 1000$  with  $T = 100$ , the average squared norm losses (scaled by 100) of the WF-SOFAR and PC factor estimates,  $\{\text{WS}, \text{PC}\}$ , are  $\{6.2, 11.6\}, \{4.6, 10.1\}, \{3.5, 9.3\}, \{2.8, 9.0\}$ , respectively. In terms of their ratios, the WF-SOFAR factor estimates become more accurate than the PC factor estimates as  $N$  rises with given  $T$ .

As expected, the accuracy of the WF-SOFAR estimator of factor loadings is uniformly superior to that of the PC estimator. This gap in accuracy becomes wider when the exponents are further from unity. For example, when  $(\alpha_1, \alpha_2) = (0.9, 0.5)$ , for  $N = 100, 200, 500, 1000$  with  $T = 500$ , the average squared norm losses (scaled by 100) for the WF-SOFAR and PC,  $\{\text{WS}, \text{PC}\}$ , are  $\{2.1, 6.9\}, \{1.7, 8.2\}, \{1.6, 11.3\}, \{1.6, 14.8\}$ , respectively. The norm loss of the WF-SOFAR stabilizes while that of the PC fast rises. However, when  $T$  grows with given  $N$ , the accuracy of both the WF-SOFAR and the PC factor loadings estimates improves (but the former is always more accurate than the latter). For example, when  $(\alpha_1, \alpha_2) = (0.9, 0.5)$ , for  $T = 100, 200, 500, 1000$  with  $N = 500$ , the average squared norm losses (scaled by 100) of the WF-SOFAR and PC factor loadings estimates,  $\{\text{WS}, \text{PC}\}$ , are  $\{10.3, 63.2\}, \{4.5, 29.5\}, \{1.6, 11.3\}, \{0.8, 5.7\}$ , respectively.

Consequently, the accuracy of the WF-SOFAR estimator of common component is uniformly superior to that of the PC estimator. For example, when  $(\alpha_1, \alpha_2) = (0.9, 0.5)$ ,  $T = N = 100, 200, 500, 1000$ , the average squared norm losses (scaled by 100) of the WF-SOFAR and PC factor component estimates,  $\{\text{WS}, \text{PC}\}$ , are  $\{19.5, 47.6\}, \{11.5, 30.0\}, \{6.2, 16.5\}, \{3.9, 10.5\}$ , respectively.

Table 3 reports the same information as Table 2, but for more challenging models with  $(\alpha_1, \alpha_2) = (0.8, 0.4), (0.5, 0.5)$  and  $(0.5, 0.4)$ . As the WF-SOFAR estimation naturally requires a larger sample size for these cases, we consider the combinations for  $N, T = 500, 1000$  only. As can be seen in the table, remarkably, even when one of the exponent is 0.4, our WF-SOFAR method provides sufficiently accurate estimates of  $\alpha_1$  and  $\alpha_2$  as well as far superior estimates of factors, factor loadings and common components to the PC method.

To summarize, the WF-SOFAR estimator performs very well when the exponents are close to unity, thus, signal of common components is high, even with a smaller sample size. When the signal of common components is weak, namely when the value(s) of exponent(s) are around 1/2 or below, the WF-SOFAR estimator is sufficiently precise in terms of norm

loss but requires a larger sample size. Significantly, even when the gap between the largest exponent and the smallest exponent is larger than the condition (10) implies (the latter must be greater than half of the former to ensure the divergence of the common components though), the WF-SOFAR estimator is sufficiently accurate in terms of norm loss and its accuracy improves as the sample size rises. Conversely, the PC estimator fails to improve the performance when  $N$  rises due to its inability to identify zero elements in sparse loadings, and consequently the PC estimator is uniformly superseded by the WF-SOFAR estimator in terms of norm loss.

## 6 Empirical Applications

In this section we provide two empirical applications. In the first subsection, the WF-SOFAR is applied to firm security returns to analyse changes in the presence of systematic risks in the market over the decades. In the second subsection we compare the forecasting performance of predictive regressions based on the factors extracted by the WF-SOFAR and WF-PC.

### 6.1 Firm security returns

In this subsection we apply our method to estimate approximate factor models using excess returns of firm securities which are used to compute the Standard & Poor's 500 (S&P 500) index of large cap U.S. equities market. In particular, we obtain the 500 securities that constitute the S&P 500 index each month over the period from January 1984 to April 2018 from Datastream. The monthly return of security  $i$  for month  $t$  is computed as  $r_{ti} = 100 \times (P_{ti} - P_{t-1,i})/P_{t-1,i} + DY_{ti}/12$ , where  $P_{ti}$  is the end-of-the-month price of the security and  $DY_{ti}$  is the per cent per annum dividend yield on the security. The one-month US treasury bill rate is chosen as the risk-free rate ( $r_{ft}$ ), which is obtained from Ken French's data library web page. The excess return is defined as  $r_{e,ti} = r_{ti} - r_{ft}$ .

Following the literature, we estimate the factor model for the standardized excess return,  $r_{e,ti}^*$ . In view of the experimental results shown earlier, we report the results based on the adaptive WF-SOFAR with BIC. For each window month,  $T = \text{Sept 1998}, \dots, \text{April 2018}$ , we chose securities that contain the data extending 120 months back ( $T = 120$ ) from  $T$ . This gives the different number of securities for each window  $T$  ( $N_T$ ). The average number of securities over the estimation windows is 443 ( $\bar{N} = 443$ ). In this exercise, we set the maximum number of factors as four. As will be shown below, three or four factors are estimated over the windows. We identify the factors and signs of the factors and factor loadings, given the estimates of the initial window month,  $T = \text{September 1989}$ , based on the correlation coefficients between the factors at  $T$  and the appropriately lagged  $T$ .<sup>11</sup>

We report the estimates of the exponents of the eigenvalues of the security return covariance matrix,  $\alpha_\ell$ ,  $\ell = 1, 2, 3, 4$ , which are associated with the four factors. Observe that, as discussed earlier, the estimated exponents are invariant to the rotation of the estimated common components. Table 4 reports the summary statistics of estimated  $\alpha_\ell$ 's and the portion of non-zero factors,  $N_{\ell T}/N_T$ ,  $\ell = 1, 2, 3, 4$ , and Figure 1 plots the estimated  $\alpha_\ell$ ,  $\ell = 1, 2, 3, 4$  over the estimation window months,  $T = \text{September 1989}, \dots, \text{April 2018}$ .

In turn we discuss the trajectories of the estimated exponents in some details by referring to Table 4 and Figure 1. The first factor does seem to be almost always "strong," in that

<sup>11</sup>For example, define  $(T - 1)$ -dimensional vector of  $\ell$ th factor of  $T$  as  $\hat{\mathbf{f}}_{\ell T} = (\hat{f}_{\ell T,1}, \hat{f}_{\ell T,2}, \dots, \hat{f}_{\ell T,T-1})'$  and that of  $T - 1$  as  $\hat{\mathbf{f}}_{\ell T-1} = (\hat{f}_{\ell T-1,2}, \hat{f}_{\ell T-1,3}, \dots, \hat{f}_{\ell T-1,T})'$ ,  $\ell = 1, \dots, r$ . For  $\hat{\mathbf{f}}_{\ell T}$ , if  $\max_{1 \leq k \leq r} |\text{corr}(\hat{\mathbf{f}}_{\ell T}, \hat{\mathbf{f}}_{kT-1})| = |\text{corr}(\hat{\mathbf{f}}_{\ell T}, \hat{\mathbf{f}}_{2T-1})|$  and  $\text{corr}(\hat{\mathbf{f}}_{\ell T}, \hat{\mathbf{f}}_{2T-1}) < 0$ , say,  $\hat{\mathbf{f}}_{2T} \equiv -\hat{\mathbf{f}}_{\ell T}$  and  $\hat{\mathbf{b}}_{i2T} \equiv -\hat{\mathbf{b}}_{i\ell T}$ .

Table 2: Performance of the WF-SOFAR (WS) and PC estimators for approximate factor models with two factor components with  $\alpha_1 \geq 0.8$  and  $\alpha_2 \geq 0.5$

Design ( $\alpha_1, \alpha_2$ )	T=100						T=200					
	(0.9, 0.9)		(0.9, 0.8)		(0.9, 0.5)		(0.9, 0.8)		(0.9, 0.5)		(0.8, 0.8)	
	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.
<b>N=100</b>												
$\hat{\alpha}_1$	0.86	0.02	0.84	0.03	0.85	0.02	0.76	0.03	0.85	0.02	0.75	0.03
$\hat{\alpha}_2$	0.85	0.02	0.78	0.04	0.54	0.07	0.75	0.03	0.52	0.07	0.52	0.07
	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC
$L^2(\hat{\mathbf{F}})_{\times 100}$	6.2	11.6	6.3	11.9	13.0	21.0	8.2	13.4	13.8	21.8	7.1	9.6
$L^2(\hat{\mathbf{B}})_{\times 100}$	9.0	9.9	8.9	12.0	9.7	31.7	9.9	15.4	10.4	38.2	5.3	8.3
$L^2(\hat{\mathbf{C}})_{\times 100}$	8.2	14.5	10.9	18.4	19.5	47.6	10.5	19.8	20.9	50.6	7.6	12.3
<b>N=200</b>												
$\hat{\alpha}_1$	0.86	0.01	0.85	0.02	0.85	0.01	0.77	0.02	0.75	0.02	0.77	0.01
$\hat{\alpha}_2$	0.86	0.01	0.77	0.03	0.54	0.05	0.76	0.02	0.52	0.05	0.77	0.02
	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC
$L^2(\hat{\mathbf{F}})_{\times 100}$	4.6	10.1	4.2	9.9	9.9	19.0	5.9	11.5	10.4	19.5	4.8	7.7
$L^2(\hat{\mathbf{B}})_{\times 100}$	9.1	10.4	8.6	12.4	9.8	41.9	10.2	17.4	10.0	50.0	5.3	9.3
$L^2(\hat{\mathbf{C}})_{\times 100}$	6.8	13.1	8.7	16.8	15.9	54.1	8.6	18.7	16.4	56.8	5.6	11.0
<b>N=500</b>												
$\hat{\alpha}_1$	0.87	0.01	0.86	0.02	0.86	0.01	0.77	0.01	0.75	0.01	0.78	0.01
$\hat{\alpha}_2$	0.86	0.01	0.77	0.02	0.55	0.04	0.76	0.01	0.52	0.04	0.77	0.01
	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC
$L^2(\hat{\mathbf{F}})_{\times 100}$	3.5	9.3	3.0	8.6	6.8	18.5	4.3	10.1	7.0	18.8	3.3	6.4
$L^2(\hat{\mathbf{B}})_{\times 100}$	9.4	11.2	9.0	13.6	10.3	63.2	10.6	20.3	10.8	74.8	5.4	10.7
$L^2(\hat{\mathbf{C}})_{\times 100}$	6.1	12.7	7.5	16.6	13.1	72.9	7.5	19.1	13.4	76.0	4.2	10.3
<b>N=1000</b>												
$\hat{\alpha}_1$	0.87	0.01	0.86	0.01	0.86	0.01	0.77	0.01	0.76	0.01	0.78	0.01
$\hat{\alpha}_2$	0.86	0.01	0.77	0.01	0.56	0.03	0.76	0.01	0.53	0.03	0.77	0.01
	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC
$L^2(\hat{\mathbf{F}})_{\times 100}$	2.8	9.0	2.2	7.8	5.1	20.0	3.4	9.4	5.2	20.1	2.4	5.8
$L^2(\hat{\mathbf{B}})_{\times 100}$	9.4	12.0	8.8	14.4	10.9	86.7	11.0	23.1	11.5	101.8	5.3	12.0
$L^2(\hat{\mathbf{C}})_{\times 100}$	6.0	12.7	6.7	16.6	12.2	96.0	7.2	20.1	12.3	99.6	3.8	10.7

Notes: The DGP is  $x_{ti} = \sum_{k=1}^r b_{ik}^0 f_{tk}^0 + \sqrt{\theta} e_{ti}$ ,  $i = 1, \dots, N$ ,  $t = 1, \dots, T$ , where  $b_{ik}^* \sim \text{i.i.d.} N(0, 1)$  for  $i = 1, \dots, N$  and  $b_{ik}^* = 0$  for  $i = N_k + 1, \dots, N$ , with  $0 < \alpha_k \leq 1$ ,  $f_{tk}^* = \rho_{fk} f_{k,t-1}^* + v_{tk}$  for  $t = 1, \dots, T$ ,  $v_{tk} \sim \text{i.i.d.} N(0, 1 - \rho_{fk}^2)$  with  $f_{0k}^* \sim \text{i.i.d.} N(0, 1)$  for  $k = 1, \dots, r$ ,  $e_{ti} = \rho_e e_{t,i-1} + \beta \varepsilon_{t,i-1} + \varepsilon_{ti}$ ,  $\varepsilon_{ti} \sim \text{i.i.d.} N(0, \sigma_{\varepsilon_{ti}}^2)$ ,  $\sigma_{\varepsilon_{ti}}^2$  is set so that  $\text{Var}(e_{ti}) = 1$ . Once  $b_{ik}^*$  and  $f_{tk}^*$  are generated, we apply Gram-Schmidt orthonormalization and form  $b_{ik}^0$  and  $f_{tk}^0$ , where  $N^{-1} \sum_{i=1}^N b_{ik}^0 b_{i\ell}^0 = 1\{k = \ell\}$  and  $T^{-1} \sum_{t=1}^T f_{tk}^0 f_{t\ell}^0 = 1\{k = \ell\}$ . We set  $r = 2$ ,  $\rho_{fk} = \rho_e = 0.5$  for  $k = 1, \dots, r$ ,  $\beta = 0.2$ ,  $\theta = 1$ . The rows  $\hat{\alpha}_1$  and  $\hat{\alpha}_2$  in the table report the averages and standard deviations (s.d.) of the WF-SOFAR estimates of  $\alpha_1$  and  $\alpha_2$ , and the rows  $L^2(\hat{\mathbf{F}})_{\times 100}$ ,  $L^2(\hat{\mathbf{B}})_{\times 100}$  and  $L^2(\hat{\mathbf{C}})_{\times 100}$  report the averages of the norm losses (multiplied by 100) of the scaled estimated factors, factor loadings and common components by the WF-SOFAR and PC over the replications. All the results are based on 1000 replications.

Table 2 (Continued):

Design ( $\alpha_1, \alpha_2$ )	T=500						T=1000					
	(0.9, 0.9)		(0.9, 0.8)		(0.9, 0.5)		(0.9, 0.8)		(0.9, 0.5)		(0.8, 0.8)	
	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.
<b>N=100</b>												
$\hat{\alpha}_1$	0.88	0.01	0.87	0.02	0.87	0.01	0.78	0.02	0.88	0.01	0.79	0.01
$\hat{\alpha}_2$	0.88	0.01	0.79	0.03	0.52	0.05	0.78	0.02	0.52	0.04	0.78	0.01
	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC
$L^2(\hat{\mathbf{F}})_{\times 100}$	4.2	5.3	4.5	5.7	11.9	13.6	6.1	7.2	11.5	12.4	5.7	6.2
$L^2(\hat{\mathbf{B}})_{\times 100}$	2.2	2.6	2.1	3.0	2.1	6.9	2.4	3.8	1.2	3.7	1.5	2.2
$L^2(\hat{\mathbf{C}})_{\times 100}$	4.1	5.5	5.9	7.5	13.4	19.2	6.1	8.0	12.7	15.6	5.5	6.5
	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.
<b>N=200</b>												
$\hat{\alpha}_1$	0.88	0.01	0.88	0.01	0.88	0.01	0.78	0.01	0.88	0.01	0.79	0.01
$\hat{\alpha}_2$	0.88	0.01	0.79	0.02	0.52	0.04	0.78	0.01	0.51	0.03	0.78	0.01
	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC
$L^2(\hat{\mathbf{F}})_{\times 100}$	2.8	4.1	2.6	3.8	8.2	9.9	4.0	5.1	8.1	8.9	3.7	4.3
$L^2(\hat{\mathbf{B}})_{\times 100}$	2.2	2.8	1.6	2.6	1.7	8.2	2.3	4.1	1.0	4.2	1.5	2.4
$L^2(\hat{\mathbf{C}})_{\times 100}$	2.6	4.0	3.5	5.3	9.5	16.8	3.9	6.1	8.9	12.5	3.4	4.5
	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.
<b>N=500</b>												
$\hat{\alpha}_1$	0.88	0.00	0.88	0.00	0.88	0.00	0.78	0.01	0.89	0.00	0.79	0.01
$\hat{\alpha}_2$	0.88	0.00	0.79	0.01	0.52	0.03	0.78	0.01	0.51	0.02	0.79	0.01
	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC
$L^2(\hat{\mathbf{F}})_{\times 100}$	1.8	3.2	1.4	2.5	5.2	7.0	2.5	3.9	5.0	5.9	2.2	2.9
$L^2(\hat{\mathbf{B}})_{\times 100}$	2.2	3.0	1.4	2.6	1.6	11.3	2.4	4.9	0.8	5.7	1.5	2.7
$L^2(\hat{\mathbf{C}})_{\times 100}$	1.7	3.2	2.0	4.1	6.2	16.5	2.4	4.9	5.6	10.7	1.9	3.2
	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.
<b>N=1000</b>												
$\hat{\alpha}_1$	0.89	0.00	0.88	0.00	0.88	0.00	0.79	0.00	0.89	0.00	0.79	0.00
$\hat{\alpha}_2$	0.88	0.00	0.79	0.01	0.53	0.02	0.78	0.00	0.52	0.02	0.79	0.00
	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC
$L^2(\hat{\mathbf{F}})_{\times 100}$	1.4	2.9	1.0	2.1	3.6	5.5	1.9	3.3	3.4	4.3	1.5	2.2
$L^2(\hat{\mathbf{B}})_{\times 100}$	2.1	3.1	1.4	2.9	1.6	14.8	2.4	5.3	0.7	7.3	1.3	2.9
$L^2(\hat{\mathbf{C}})_{\times 100}$	1.4	2.9	1.6	3.8	4.6	18.1	1.9	4.7	3.9	10.5	1.3	2.7
	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.



Table 3: Performance of the WF-SOFAR (WS) and PC estimators for approximate factor models with two factor components with  $(\alpha_1, \alpha_2) = (0.8, 0.4), (0.5, 0.5), (0.5, 0.4)$

Design $(\alpha_1, \alpha_2)$	<b>T=500</b>						<b>T=1000</b>					
	(0.8, 0.4)		(0.5, 0.5)		(0.5, 0.4)		(0.8, 0.4)		(0.5, 0.5)		(0.5, 0.4)	
	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.	mean	s.d.
<b>N=500</b>												
$\hat{\alpha}_1$	0.78	0.01	0.48	0.02	0.47	0.03	0.78	0.01	0.48	0.02	0.47	0.03
$\hat{\alpha}_2$	0.44	0.04	0.48	0.02	0.41	0.04	0.43	0.03	0.48	0.02	0.40	0.04
	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC
$L^2(\hat{\mathbf{F}})_{\times 100}$	9.3	13.4	10.5	12.9	13.4	17.9	9.2	11.1	10.0	11.0	13.1	15.2
$L^2(\hat{\mathbf{B}})_{\times 100}$	2.0	25.9	3.2	30.6	4.6	48.3	1.1	12.8	1.9	15.4	2.9	24.4
$L^2(\hat{\mathbf{C}})_{\times 100}$	10.9	34.3	10.6	30.0	17.3	48.6	10.0	21.3	10.0	19.4	16.0	31.1
<b>N=1000</b>												
$\hat{\alpha}_1$	0.78	0.00	0.49	0.01	0.48	0.02	0.79	0.00	0.49	0.01	0.48	0.02
$\hat{\alpha}_2$	0.44	0.03	0.48	0.02	0.40	0.03	0.43	0.03	0.49	0.01	0.40	0.03
	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC	WS	PC
$L^2(\hat{\mathbf{F}})_{\times 100}$	6.8	11.7	7.5	10.1	9.7	15.2	6.7	8.9	7.2	8.3	9.5	12.0
$L^2(\hat{\mathbf{B}})_{\times 100}$	2.0	36.1	3.1	40.4	3.7	65.6	1.0	17.4	1.8	20.0	2.3	32.2
$L^2(\hat{\mathbf{C}})_{\times 100}$	8.3	41.5	7.8	33.3	13.0	57.4	7.5	23.1	7.0	19.4	12.0	32.9

Notes: See the notes to Table 2.

the divergence rate  $N_1$  is very close to  $N$ . As reported in Table 4, the average of  $\alpha_1$  over the month windows is 0.995 and standard deviation is very small (0.004) with the minimum value of 0.979. Actually, the values of the factor loadings to this factor have the same sign, which strongly suggests that this is the market factor. Apart from the first factor, which is always strong, the strengths of the common components vary over the months and can become quite weak. For example, for the second to the fourth factors, the maximum portion of nonzero factor loadings is between 44.5% and 59.5%, whilst the minimum portion is merely 12.0% to 17.6%. Furthermore, the divergence rates are very different over the factors. For example, for the window month of March 1998,  $\{\hat{\alpha}_1, \hat{\alpha}_2, \hat{\alpha}_3\} = \{0.991, 0.774, 0.653\}$ , and the corresponding numbers of non-zero factors are 425, 113 and 54 out of 450 securities. These strongly imply a potentially substantial efficiency gain in estimation of the approximate factor models through our WF-SOFAR method over the PC method.

In line with the well-observed phenomenon that the correlation among the securities in the financial market rises during periods of turmoil, sharp rises of exponents in some months can be observed. For example,  $\alpha_2$  goes up sharply around February 2000 then rises gradually. This period corresponds to the peak of the dot-com bubble and its burst on March 2000 (the main contributor to the factor loadings of the second factor is Technology industry, see Appendix D). Similarly, a sharp rise of  $\alpha_3$  is observed from July 2008 to April 2009. This period coincides with the 2008 financial crisis. In just ten months, it goes up by 0.12, from 0.74 to 0.86 (one of the main contributors to the factor loadings of the third factor is the Financial industry, see Appendix D).

It is also interesting that the orders in terms of values of the exponents,  $\alpha_2$ ,  $\alpha_3$ , and  $\alpha_4$ , change over the period. In particular, from September 1989,  $\alpha_2$  is larger than  $\alpha_3$  most of the time until December 2010, then  $\alpha_3$  is almost always larger than  $\alpha_2$ . Since the estimate of  $\alpha_4$  first appeared in February 2004, it was mostly smaller than other exponents. It is estimated every month from March 2010 onward, seemingly becoming more and more strong toward the latest month, April 2018. After the sharp one-off drop in February 2015,<sup>12</sup>  $\alpha_4$  rises to become the highest next to the first factor from November 2016 onward.

Table 4: Summary statistics of the estimated exponents of factor loadings,  $\alpha_{\ell T}$ , and the portion of non-zero factor loadings,  $N_{\ell T}/N_T$ ,  $\ell = 1, 2, 3, 4$ , from September 1989 to April 2018.

	Exponents of Loadings				Portion of Non-zero Loadings			
	$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_4$	$N_1/N$	$N_2/N$	$N_3/N$	$N_4/N$
<b>mean</b>	0.995	0.824	0.770	0.781	97.1%	36.2%	26.2%	27.3%
<b>s.d.</b>	0.004	0.046	0.045	0.028	2.4%	9.6%	7.0%	5.0%
<b>max</b>	1.000	0.895	0.860	0.854	100.0%	59.5%	44.5%	49.7%
<b>min</b>	0.979	0.713	0.653	0.665	88.2%	17.6%	12.0%	13.1%

Notes: The estimated  $\alpha_\ell$ ,  $\ell = 1, 2, 3, 4$  for the each month of 120 months window,  $T =$  September 1989,...,April 2018.

## 6.2 Forecasting bond yields

In this subsection we consider out-of-sample performance of forecasting regressions for bond yields using extracted factors via our WF-SOFAR method and the PC method, from a large

<sup>12</sup>This coincides with the period at bottom of the biggest sharp fall of oil price between 2014-2015.

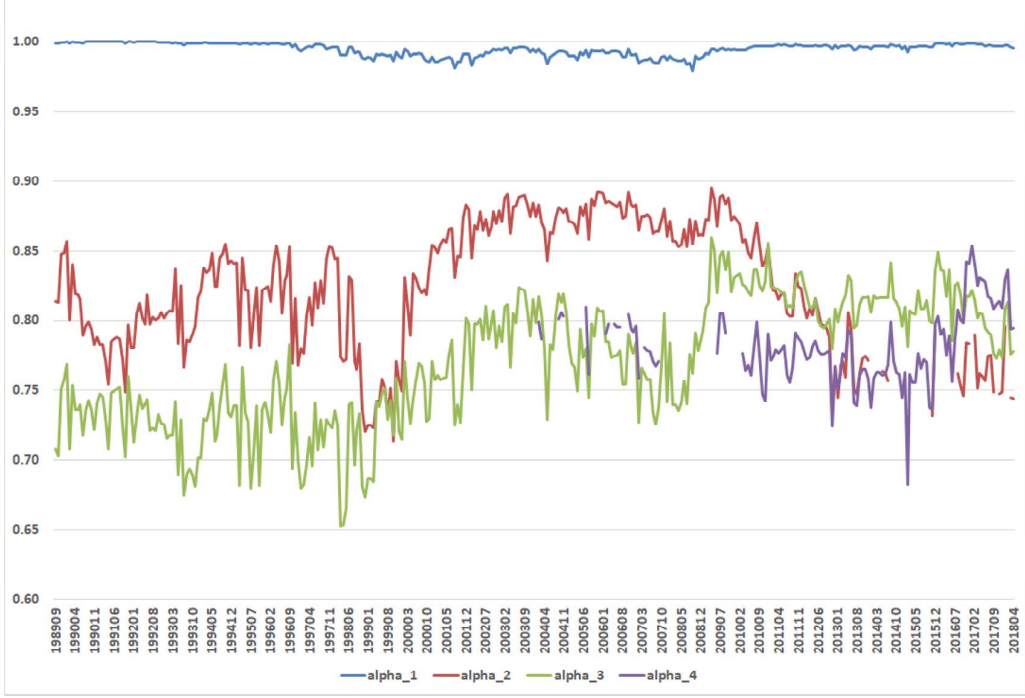


Figure 1: Plot of the estimated  $\alpha_k$ 's from September 1989 to April 2018. The estimated  $\alpha_\ell$ ,  $\ell = 1, 2, 3, 4$  for the each month of 120 months window,  $T = \text{September 1989}, \dots, \text{April 2018}$ .

number of macroeconomic (prediction) variables in line with the analysis of Ludvigson and Ng (2009). We use the same data set which is provided by Ludvigson and Ng.<sup>13</sup> Specifically, the data consists of the continuously compounded (log) annual excess returns on an  $n$ -year discount bond at month  $t$ ,  $y_t^{(n)}$ , and a balanced panel of  $i = 1, \dots, 132$  monthly macroeconomic series at month  $t$ ,  $x_{ti}$ , spanning the period from January 1964 to December 2003. We consider the maturities  $n = 2, 3, 4, 5$ .<sup>14</sup>

We consider one-year-ahead out of sample forecast comparisons of continuously compounded annual log excess bond returns,  $y_{t+12}^{(n)}$ . In order to minimize possible adverse effects of structural breaks, we set the rolling window at 252 months. The forecast comparison procedure is explained below. For the  $T$ th month rolling window and maturity  $n$ , we extract factors  $\{\hat{f}_{tk}\}_{k=1}^{\hat{r}_T}$  from  $x_{ti}$  via our WF-SOFAR and the PC,  $i = 1, \dots, N = 132$ ,  $t = T, \dots, T_T - 12$ , where  $t$  denotes months from January 1964 to December 2003,  $T$  and  $T_T$  denote the start and end months of the  $T$ th rolling window, respectively. Observe that the number of factors is estimated for each estimation window to avoid using “future” information.<sup>15</sup> Then, run the predictive regression

$$y_{t+12}^{(n)} = \tilde{\beta}_0^{(n)} + \sum_{k=1}^{\hat{r}_T} \tilde{\beta}_k^{(n)} \hat{f}_{tk} + \tilde{\varepsilon}_t^{(n)}, \quad t = T, \dots, T_T - 12, \quad n = 2, 3, 4, 5$$

<sup>13</sup>The data file is obtained from Sydney Ludvigson's web page: <https://www.sydneyludvigson.com/s/RFS2009-u1e1.xls>

<sup>14</sup>For more details of the data, see Section 3 of Ludvigson and Ng (2009).

<sup>15</sup>In another experiment, we estimated the number of factors using a whole sample period and implemented a similar exercise. The forecast based on our estimator uniformly outperformed that based on the PC estimator.

and obtain the forecast error

$$\hat{\varepsilon}_{T_{\tau}+12|T_{\tau}}^{(n)} = y_{T_{\tau}+12}^{(n)} - \hat{y}_{T_{\tau}+12|T_{\tau}}^{(n)},$$

with  $\hat{y}_{T_{\tau}+12|T_{\tau}}^{(n)} = \tilde{\beta}_0^{(n)} + \sum_{k=1}^{\hat{r}_{T_{\tau}}} \tilde{\beta}_k^{(n)} \hat{f}_{T_{\tau}k}$ . This produces  $H = 217$  forecast errors. To estimate the number of factors for each window, we set the maximum number of factors to nine and use the ED estimator. The estimated number of factors varies from one to six over the forecast windows. In Table 1, we report the mean absolute deviation of the forecast errors,  $MAE^{(n)} = H^{-1} \sum_{s=1}^H |\hat{\varepsilon}_{s|s-1}^{(n)}|$  for  $n = 2, 3, 4, 5$ , for our WF-SOFAR and the PC, and Diebold-Mariano forecasting performance test statistics with associated p-values, based on the MAEs.<sup>16</sup> As can be seen, the MAEs of the WF-SF are smaller than those of the PC for all the maturities. The Diebold-Mariano forecasting performance test strongly rejects the null of the same forecasting performance for all the maturities, in favour of the alternative that our method outperforms the PC method. The average values of exponents over the windows are  $\{\alpha_1, \alpha_2, \alpha_3, \alpha_4, \alpha_5, \alpha_6\} = \{0.92, 0.82, 0.87, 0.78, 0.77, 0.74\}$ , which suggests that even the (first) strongest factor component is not strictly strong. As is evidenced in the previous section, the accuracy of our estimator is much higher than the PC estimator under such situations, and the better forecasting performance may not be too surprising in this empirical exercise.

Table 5: Mean absolute forecast errors and Diebold-Mariano forecast comparison test result

	WS	PC	Diebold-Mariano Statistic	[p-value]
$y_{t+12}^{(2)}$	1.164	1.191	-3.58	[0.0003]
$y_{t+12}^{(3)}$	2.304	2.354	-3.54	[0.0004]
$y_{t+12}^{(4)}$	3.354	3.429	-3.73	[0.0002]
$y_{t+12}^{(5)}$	4.197	4.278	-3.20	[0.0014]

Notes: For the computation of the long-run variance for the Diebold-Mariano test statistic of Diebold and Mariano (1995), the window is chosen by the Schwert criterion with the maximum lag of 14.

## 7 Conclusion

This paper considers the determination of number of factors and efficient estimation of the large dimensional approximate factor models of Chamberlain and Rothschild (1983)), called weak factor (WF) models, in which the  $r$  largest eigenvalues of a data covariance matrix grow as  $N$  rises with possibly different diverging rates  $N^{\alpha_k}$ ,  $0 < \alpha_k \leq 1$  for  $k = 1, \dots, r$ . We consider the WF structure induced by *sparse* factor loadings  $\mathbf{B}^0$  that leads to  $\lambda_k(\mathbf{B}^{0'}\mathbf{B}^0) \asymp N^{\alpha_k}$ . The proposed *weak factor sparse orthogonal factor regression* (WF-SOFAR) estimator and its adaptive version enable us to consistently estimate the WF models. As theoretical contributions, the non-asymptotic error bound and the rate of convergence for the WF-SOFAR estimators are derived. For the adaptive WF-SOFAR estimator, we have established the factor selection consistency and consistent estimation of each exponent  $\alpha_k$  of the  $r$  divergence rates. The proposed estimator is expected to be more efficient than

<sup>16</sup>We computed the Diebold-Mariano test statistic based on mean squared errors, and the null was more strongly rejected.

the widely used principal component (PC) estimator. In line with our theoretical results, the Monte Carlo experiments show that the ED estimator of Onatski (2010) for the number of factors  $r$  is reliable while other widely used estimators such as Bai and Ng (2002) and Ahn and Horenstein (2013) can be unreliable. Furthermore, the evidence demonstrates that the WF-SOFAR estimator uniformly dominates the PC estimator in terms of the norm loss, including under the experimental design most favorable to the PC estimator, in which the exponents of the divergence rates are 0.9.

Recently estimation of a hierarchical factor structure or a multi-level factor structure has been gaining serious interest in the literature. Ando and Bai (2017) and Choi, Kim, Kim, and Kwark (Choi et al.) consider factor models with two types of factors, global factors and local factors. The factor loadings of global factors are non-zero values for all the cross-section units, whereas the local factors have non-zero loadings among the cross-section units of specific cross sectional groups. Ando and Bai (2017) and Choi, Kim, Kim, and Kwark (Choi et al.) propose sequential procedures to identify the global and local factors separately. In fact, the WF structure nests the hierarchical factor structure and hence our WF-SOFAR method can be applied to readily estimate such models. In contrast to existing approaches, given the total number of global and local factors, our approach permits us to consistently estimate the number of local groups, the number of global and local factors and its memberships in one go. For further information and additional simulation results, see Appendix C.

Having provided the consistency result of the WF-SOFAR estimator in this paper, the statistical inference for the high-dimensional WF models is an important research agenda. The application of the method named *debiasing* (desparsification) by Javanmard and Montanari (2014) and van de Geer et al. (2014) to the WF-SOFAR estimator seems the most promising pathway. This methodology will correct the Lasso estimator to remove the bias based on the Karush-Kuhn-Tucker (KKT) conditions. However, it is nontrivial to establish how the WF-SOFAR estimator can be debiased, and this is our future challenge.

In this paper we have focused on the estimation of the common factors and the exponents of the divergence rates of the  $r$  largest eigenvalues. It is of interest to estimate the stock return covariance matrix for optimal portfolio allocation and portfolio risk assessment. This can be achieved by consistently estimating the covariance matrix of idiosyncratic errors, in line with Fan et al. (2008) and Fan et al. (2011), which is an interesting extension of this paper.

Another possible extension of interest is to consider the estimation of panel data models with unobservable multiple interactive effects. Pesaran (2006) and Bai (2009), among others, develop the estimation methods of the panel data model:

$$y_{ti} = \mathbf{x}_{ti}'\boldsymbol{\beta} + u_{ti}, \quad u_{ti} = \mathbf{f}_t'\mathbf{b}_i + \varepsilon_{ti}.$$

For the PC based estimators, such as Bai (2009),  $u_{it}$  is typically assumed to have the strong factor structure (i.e.,  $\sum_{i=1}^N \mathbf{b}_i \mathbf{b}_i' / N$  tends to a fixed matrix), which may not hold in practice, and the WF structure seems more appropriate. The iterative procedure proposed by Bai (2009) based on the WF-SOFAR estimation of  $\mathbf{f}_t'\mathbf{b}_i$ , instead of the PC estimation, would potentially improve the precision of the estimates of  $\boldsymbol{\beta}$ .

## References

- Ahn, S. C. and A. R. Horenstein (2013). Eigenvalue ratio test for the number of factors. *Econometrica* 81, 1203–1227.

- Amengual, D. and M. Watson (2007). Consistent estimation of the number of dynamic factors in a large  $N$  and  $T$  panel. *Journal of Business & Economic Statistics* 25, 91–96.
- Ando, T. and J. Bai (2017). Clustering huge number of financial time series: A panel data approach with high-dimensional predictors and factor structures. *Journal of the American Statistical Association* 112, 1182–1198.
- Bai, J. (2003). Inferential theory for factor models of large dimensions. *Econometrica* 71, 135–171.
- Bai, J. (2009). Panel data models with interactive fixed effects. *Econometrica* 77, 1229–1279.
- Bai, J. and S. Ng (2002). Determining the number of factors in approximate factor models. *Econometrica* 70, 191–221.
- Bai, J. and S. Ng (2006). Confidence intervals for diffusion index forecasts and inference with factor-augmented regressions. *Econometrica* 74, 1133–1150.
- Bai, J. and S. Ng (2013). Principal components estimation and identification of static factors. *Journal of Econometrics* 176, 18–29.
- Bai, J. and S. Ng (2017). Principal components and regularized estimation of factor models. *arXiv:1708.08137v2*.
- Bailey, N., G. Kapetanios, and M. H. Pesaran (2016). Exponent of cross-sectional dependence: Estimation and inference. *Journal of Applied Econometrics* 31, 929–960.
- Bryzgalova, S. (2016). Spurious factors in linear asset pricing models. *mimeo*.
- Candès, E. J., X. Li, Y. Ma, and J. Wright (2011). Robust principal component analysis? *Journal of the ACM* 58, 11:1–11:37.
- Caner, M. and X. Han (2014). Selecting the correct number of factors in approximate factor models: The large panel case with group bridge estimators. *Journal of Business & Economic Statistics* 32, 359–374.
- Chamberlain, G. and M. Rothschild (1983). Arbitrage, factor structure and mean-variance analysis in large asset markets. *Econometrica* 51, 1281–1304.
- Choi, I., D. Kim, Y. J. Kim, and N.-S. Kwark. A multilevel factor model: Identification, asymptotic theory and applications. *Journal of Applied Econometrics* 33, 355–377.
- Connor, G. and R. A. Korajczyk (1986). Performance measurement with the arbitrage pricing theory: A new framework for analysis. *Journal of Financial Economics* 15, 373–394.
- Connor, G. and R. A. Korajczyk (1993). A test for the number of factors in an approximate factor model. *Journal of Finance* 48, 1263–1291.
- De Mol, C., D. Giannone, and L. Reichlin (2008). Forecasting using a large number of predictors: Is Bayesian shrinkage a valid alternative to principal components? *Journal of Econometrics* 146, 318–328.
- Diebold, F. X. and R. S. Mariano (1995). Comparing predictive accuracy. *Journal of Business & Economic Statistics* 13, 253–263.

- Fama, E. F. and K. R. French (2015). A five-factor asset pricing model. *Journal of Financial Economics* 116, 1–22.
- Fan, J., Y. Fan, and E. Barut (2014). Adaptive robust variable selection. *Annals of Statistics* 42, 324–351.
- Fan, J., Y. Fan, and J. Lv (2008). High dimensional covariance matrix estimation using a factor model. *Journal of Econometrics* 147, 186–197.
- Fan, J. and R. Li (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association* 96, 1348–1360.
- Fan, J., Y. Liao, and M. Mincheva (2011). High-dimensional covariance matrix estimation in approximate factor models. *Annals of Statistics* 39, 3320–3356.
- Fan, J., Y. Liao, and M. Mincheva (2013). Large covariance estimation by thresholding principal orthogonal complements. *Journal of the Royal Statistical Society Series B* 75, 603–680.
- Fan, J. and J. Lv (2011). Nonconcave penalized likelihood with NP-dimensionality. *IEEE Transactions on Information Theory* 57, 5467–5484.
- Fan, J., K. Wang, Y. Zhong, and Z. Zhu (2018). Robust high-dimensional factor models with applications to statistical machine learning. *arXiv:1808.03889v1*.
- Freyaldenhoven, S. (2018). A generalized factor model with local factors. *mimeo*.
- Hallin, M. and R. Liška (2007). Determining the number of factors in the general dynamic factor model. *Journal of the American Statistical Association* 102, 603–617.
- Javanmard, A. and A. Montanari (2014). Confidence intervals and hypothesis testing for high-dimensional regression. *Journal of Machine Learning Research* 15, 2869–2909.
- Johnstone, I. M. and A. Y. Lu (2009). On consistency and sparsity for principal components analysis in high dimensions. *Journal of the American Statistical Association* 104, 682–693.
- Kapetanios, G. (2010). A testing procedure for determining the number of factors in approximate factor models with large datasets. *Journal of Business & Economic Statistics* 28, 397–409.
- Kock, A. B. and H. Tang (2019). Uniform inference in high-dimensional dynamic panel data models with approximately sparse fixed effects. *Econometric Theory* 35, 295–359.
- Leeb, H. and B. M. Pötscher (2005). Model selection and inference: Facts and fiction. *Econometric Theory* 21, 21–59.
- Leeb, H. and B. M. Pötscher (2006). Can one estimate the conditional distribution of post-model-selection estimators? *Annals of Statistics* 34, 2554–2591.
- Leeb, H. and B. M. Pötscher (2008). Sparse estimators and the oracle property, or the return of hedges’ estimator. *Journal of Econometrics* 142, 201–211.
- Lettau, M. and M. Pelger (2018). Estimating latent asset-pricing factors. *NBER Working Paper 24618*.

- Ludvigson, C. S. and S. Ng (2009). Macro factors in bond risk premia. *Review of Financial Studies* 22, 5027–5067.
- Onatski, A. (2010). Determining the number of factors from empirical distribution of eigenvalues. *Review of Economics and Statistics* 92, 1004–1016.
- Onatski, A. (2012). Asymptotics of the principal components estimator of large factor models with weakly influential factors. *Journal of Econometrics* 168, 244–258.
- Pesaran, H. (2006). Estimation and inference in large heterogeneous panels with a multifactor error structure. *Econometrica* 74, 967–1012.
- Pötscher, B. M. and H. Leeb (2009). On the distribution of penalized maximum likelihood estimators: The LASSO, SCAD, and thresholding. *Journal of Multivariate Analysis* 100, 2065–2082.
- Rigollet, P. and J.-C. Hütter (2017). *High Dimensional Statistics*. Massachusetts Institute of Technology, MIT Open CourseWare.
- Rudelson, M. and R. Vershynin (2013). Hanson-wright inequality and sub-gaussian concentration. *Electronic Communications in Probability* 18, 1–9.
- Stock, J. H. and M. W. Watson (2002). Forecasting using principal components from a large number of predictors. *Journal of the American Statistical Association* 97, 1167–1179.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B*, 267–288.
- Uematsu, Y., Y. Fan, K. Chen, J. Lv, and W. Lin (2019). SOFAR: large-scale association network learning. *IEEE Transactions on Information Theory*, to appear.
- Uematsu, Y. and S. Tanaka (2019). High-dimensional macroeconomic forecasting and variable selection via penalized regression. *Econometrics Journal* 22, 34–56.
- van de Geer, S., P. Bühlmann, Y. Ritov, and R. Dezeure (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *Annals of Statistics* 42, 1166–1202.
- Vershynin, R. (2012). Introduction to the non-asymptotic analysis of random matrices. In Y. C. Eldar and G. Kutyniok (Eds.), *Compressed Sensing: Theory and Practice*, pp. 210–268. Cambridge University Press.
- Zhang, C.-H. (2014). Nearly unbiased variable selection under minimax concave penalty. *Annals of Statistics* 38, 894–942.
- Zhang, C.-H. and S. S. Zhang (2014). Confidence intervals for low dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society Series B* 76, 217–242.



## Appendix

### A Proofs of the Main Results

*Proof of Lemma 1.* (a) The  $t$ th row of  $\mathbf{E}$ ,  $\mathbf{e}_t \in \mathbb{R}^N$ , is specified as  $\mathbf{e}_t = \sum_{\ell=0}^L \Phi_\ell \boldsymbol{\varepsilon}_{t-\ell}$ , where  $\boldsymbol{\varepsilon}_t \in \mathbb{R}^N$  is composed of i.i.d.  $\text{subG}(\sigma_\varepsilon^2)$  by Assumption 3. We also define  $\tilde{\mathbf{E}}_\ell = (\boldsymbol{\varepsilon}_{1-\ell}, \dots, \boldsymbol{\varepsilon}_{T-\ell})' \in \mathbb{R}^{T \times N}$ . Then, we can write  $\mathbf{E} = \sum_{\ell=0}^{L_n} \tilde{\mathbf{E}}_\ell \Phi_\ell'$ , so that the spectral norm is bounded as

$$\|\mathbf{E}\|_2 \leq \sum_{\ell=0}^{L_n} \|\tilde{\mathbf{E}}_\ell\|_2 \|\Phi_\ell\|_2 \leq \max_{\ell \in \{0, \dots, L_n\}} \|\tilde{\mathbf{E}}_\ell\|_2 \sum_{\ell=0}^{\infty} \|\Phi_\ell\|_2.$$

By Assumption 3, the last infinite sum is bounded from above. Because of the union bound and sub-Gaussianity (see Section 4 and Theorem 5.39 of Vershynin 2012), there is a positive constant  $M$  such that

$$\begin{aligned} & \mathbb{P} \left( \max_{\ell \in \{0, \dots, L_n\}} \left\| (N \vee T)^{-1/2} \tilde{\mathbf{E}}_\ell \right\|_2 > M \right) \\ & \leq L_n \max_{\ell \in \{0, \dots, L_n\}} \mathbb{P} \left( \left\| (N \vee T)^{-1/2} \tilde{\mathbf{E}}_\ell \right\|_2 > M \right) \\ & \leq 2(N \vee T)^\nu \exp(-|O(N \vee T)|) = \exp(-|O(N \vee T)|), \end{aligned}$$

where the last inequality holds since  $\nu$  is a fixed positive constant. Thus,  $\|(N \vee T)^{-1/2} \mathbf{E}\|_2$  is bounded by a constant with probability at least  $1 - \exp(-|O(N \vee T)|)$ .

(b) By the definition, the  $(t, k)$ th element of  $\mathbf{E}\mathbf{B}^0$  is given by  $\mathbf{e}_t' \mathbf{b}_k^0 = \sum_{\ell=0}^{L_n} \boldsymbol{\varepsilon}_{t-\ell}' \Phi_\ell' \mathbf{b}_k^0$ . Let  $\tilde{b}_{\ell k, i}$  denote the  $i$ th element of  $\Phi_\ell' \mathbf{b}_k^0$ . Then, we have

$$\begin{aligned} \|\mathbf{E}\mathbf{B}^0\|_{\max} &= \max_{t \in \{1, \dots, T\}, k \in \{1, \dots, r\}} \left| \sum_{\ell=0}^{L_n} \sum_{i=1}^N \varepsilon_{t-\ell, i} \tilde{b}_{\ell k, i} \right| \\ &\leq \sum_{\ell=0}^{L_n} \max_{t \in \{1, \dots, T\}, k \in \{1, \dots, r\}} \left| \|\Phi_\ell' \mathbf{b}_k\|_2^{-1} \sum_{i=1}^N \varepsilon_{t-\ell, i} \tilde{b}_{\ell k, i} \right| \|\Phi_\ell' \mathbf{b}_k\|_2 \\ &\leq \max_{t \in \{1, \dots, T\}, k \in \{1, \dots, r\}, \ell \in \{0, \dots, L_n\}} \left| \|\Phi_\ell' \mathbf{b}_k\|_2^{-1} \sum_{i=1}^N \varepsilon_{t-\ell, i} \tilde{b}_{\ell k, i} \right| \max_k \sum_{\ell=0}^{L_n} \|\Phi_\ell\|_2 \|\mathbf{b}_k\|_2 \\ &\lesssim N_1^{1/2} \max_{t, k, \ell} \left| \|\Phi_\ell' \mathbf{b}_k\|_2^{-1} \sum_{i=1}^N \varepsilon_{t-\ell, i} \tilde{b}_{\ell k, i} \right| \sum_{\ell=0}^{\infty} \|\Phi_\ell\|_2. \end{aligned}$$

Since  $\{\varepsilon_{t-\ell, i} \tilde{b}_{\ell k, i}\}_{i=1}^N$  is a sequence of i.i.d.  $\text{subG}(\sigma_\varepsilon^2 \tilde{b}_{\ell k, i}^2)$  for each  $t, k, \ell$ , we can further see that  $\|\Phi_\ell' \mathbf{b}_k\|_2^{-1} \sum_{i=1}^N \varepsilon_{t-\ell, i} \tilde{b}_{\ell k, i} \sim \text{subG}(\sigma_\varepsilon^2)$  by Lemma 2(b). Thus, the union bound yields

$$\begin{aligned} & \mathbb{P} \left( \max_{t, k, \ell} \left| \|\Phi_\ell' \mathbf{b}_k\|_2^{-1} \sum_{i=1}^N \varepsilon_{t-\ell, i} \tilde{b}_{\ell k, i} \right| > x \right) \\ & \leq rT(L_n + 1) \max_{t, k, \ell} \mathbb{P} \left( \left| \|\Phi_\ell' \mathbf{b}_k\|_2^{-1} \sum_{i=1}^N \varepsilon_{t-\ell, i} \tilde{b}_{\ell k, i} \right| > x \right) \\ & \leq 2r(N \vee T)^{\nu+1} \exp \left( -\frac{x^2}{2\sigma_\varepsilon^2} \right). \end{aligned}$$

Setting  $x = (2\sigma_\varepsilon^2(2\nu + 1)\log(N \vee T))^{1/2}$  leads to

$$\max_{t,k,\ell} \left\| \Phi'_\ell \mathbf{b}_k \right\|_2^{-1} \sum_{i=1}^N \varepsilon_{t-\ell,i} \tilde{b}_{\ell k,i} \Big| \leq (2\sigma_\varepsilon^2(2\nu + 1)\log(N \vee T))^{1/2},$$

which holds with probability at least  $1 - O((N \vee T)^{-\nu})$ . This together with the first inequality achieves the result.

(c) Let  $\tilde{\mathbf{Z}}_\ell = (\zeta_{1-\ell}, \dots, \zeta_{T-\ell})' \in \mathbb{R}^{T \times r}$ . Then, by Assumptions 1 and 3, we can write  $\mathbf{E}'\mathbf{F} = \sum_{\ell,m=0}^L \Phi_\ell \tilde{\mathbf{E}}'_\ell \tilde{\mathbf{Z}}_m \Psi'_m$ . By the triangle inequality and property of matrix norms, we observe that

$$\begin{aligned} \|\mathbf{E}'\mathbf{F}\|_{\max} &\leq \sum_{\ell,m=0}^{L_n} \|\Phi_\ell \tilde{\mathbf{E}}'_\ell \tilde{\mathbf{Z}}_m \Psi'_m\|_{\max} \\ &\leq r^{1/2} \sum_{\ell,m=0}^{L_n} \|\Psi_m\|_2 \max_{i \in \{1, \dots, N\}, k \in \{1, \dots, r\}} \left| \phi'_{\ell,i} \tilde{\mathbf{E}}'_\ell \zeta_{m,k} \right| \\ &\leq r^{1/2} \sum_{\ell,m=0}^{L_n} \|\Psi_m\|_2 \max_{i,k} \left\| \phi_{\ell,i} \right\|_2^{-1} \left| \phi'_{\ell,i} \tilde{\mathbf{E}}'_\ell \zeta_{m,k} \right| \max_i \|\phi_{\ell,i}\|_2 \\ &\leq r^{1/2} \max_{\ell,m,i,k} \left\| \phi_{\ell,i} \right\|_2^{-1} \left| \phi'_{\ell,i} \tilde{\mathbf{E}}'_\ell \zeta_{m,k} \right| \sum_{\ell,m=0}^{L_n} \|\Psi_m\|_2 \max_i \|\phi_{\ell,i}\|_2 \\ &\leq r^{1/2} \max_{\ell,m,i,k} \left\| \phi_{\ell,i} \right\|_2^{-1} \left| \phi'_{\ell,i} \tilde{\mathbf{E}}'_\ell \zeta_{m,k} \right| \sum_{m=0}^{\infty} \|\Psi_m\|_2 \sum_{\ell=0}^{\infty} \|\Phi_\ell\|_2, \end{aligned}$$

where  $\phi'_{\ell,i}$  and  $\zeta_{m,k}$  are the  $i$ th row vector of  $\Phi_\ell$  and  $k$ th column vector of  $\tilde{\mathbf{Z}}_m$ , respectively. We can see that for each  $i$  and  $\ell$ , the row vector

$$\phi'_{\ell,i} \tilde{\mathbf{E}}'_\ell = \left( \sum_{j=1}^N \phi_{\ell,i,j} \varepsilon_{1-\ell,j}, \dots, \sum_{j=1}^N \phi_{\ell,i,j} \varepsilon_{T-\ell,j} \right)$$

is composed of i.i.d.  $\text{subG}(\sigma_\varepsilon^2 \|\phi_{\ell,i}\|_2^2)$ . Since  $\zeta_{m,k} = (\zeta_{1-m,k}, \dots, \zeta_{T-m,k})'$  consists of i.i.d.  $\text{subG}(\sigma_\zeta^2)$ , Lemma 2(a) entails that

$$\left\| \phi_{\ell,i} \right\|_2^{-1} \phi'_{\ell,i} \tilde{\mathbf{E}}'_\ell \zeta_{m,k} = \sum_{t=1}^T \left( \left\| \phi_{\ell,i} \right\|_2^{-1} \sum_{j=1}^N \phi_{\ell,i,j} \varepsilon_{t-\ell,j} \right) \zeta_{t-m,k}$$

is the sum of i.i.d.  $\text{subE}(4e\sigma_\varepsilon\sigma_\zeta)$ . Therefore, the union bound and Bernstein's inequality for the sum of sub-exponential random variables give

$$\begin{aligned} &\mathbb{P} \left( \max_{\ell,m,i,k} \left| T^{-1} \left\| \phi_{\ell,i} \right\|_2^{-1} \phi'_{\ell,i} \tilde{\mathbf{E}}'_\ell \zeta_{m,k} \right| > x \right) \\ &\leq rN(L_n + 1)^2 \max_{\ell,m,i,k} \mathbb{P} \left( \left| T^{-1} \left\| \phi_{\ell,i} \right\|_2^{-1} \phi'_{\ell,i} \tilde{\mathbf{E}}'_\ell \zeta_{m,k} \right| > x \right) \\ &\leq 2r(N \vee T)^{2\nu+1} \exp \left\{ -\frac{T}{2} \left( \frac{x^2}{16e^2\sigma_\varepsilon^2\sigma_\zeta^2} \wedge \frac{x}{4e\sigma_\varepsilon\sigma_\zeta} \right) \right\} \end{aligned}$$

for all  $x > 0$ . Putting  $x = \left(32e^2\sigma_\varepsilon^2\sigma_\zeta^2(3\nu+1)T^{-1}\log(N \vee T)\right)^{1/2}$  gives

$$\max_{\ell,m,i,k} \left| \|\phi_{\ell,i}\|_2^{-1} \phi'_{\ell,i} \tilde{\mathbf{E}}'_\ell \zeta_{m,k} \right| \leq \left(32e^2\sigma_\varepsilon^2\sigma_\zeta^2(3\nu+1)T\log(N \vee T)\right)^{1/2},$$

which holds with probability at least  $1 - O((N \vee T)^{-\nu})$ . Combining this with the first bound yields the result.

(d) To obtain the result, we apply the Hanson–Wright inequality in Rudelson and Vershynin (2013). Let  $\boldsymbol{\xi} = (\xi_1, \dots, \xi_m)' \in \mathbb{R}^m$  denote a random vector of  $m$  independent copies of  $\varepsilon \sim \text{subG}(\sigma_\varepsilon^2)$ . Then the inequality states that for any (nonrandom) matrix  $\mathbf{M} \in \mathbb{R}^{m \times m}$ ,

$$\mathbb{P}(|\boldsymbol{\xi}'\mathbf{M}\boldsymbol{\xi} - \mathbb{E}\boldsymbol{\xi}'\mathbf{M}\boldsymbol{\xi}| > u) \leq 2 \exp \left\{ -c \min \left( \frac{u^2}{K^4 \|\mathbf{M}\|_{\text{F}}^2}, \frac{u}{K^2 \|\mathbf{M}\|_2} \right) \right\}, \quad (\text{A.1})$$

where  $c$  and  $K$  are positive constants such that  $\sup_{k \geq 1} k^{-1/2} (\mathbb{E}|\varepsilon|^k)^{1/k} \leq K$ . In our setting, we can take  $K = 3\sigma_\varepsilon^2$  (e.g., Rigollet and Hütter (2017), Lemma 1.4).

Let  $\phi'_{\ell,i}$  denote the  $i$ th row vector of  $\boldsymbol{\Phi}_\ell$ . Then we have

$$\begin{aligned} \max_i \left| T^{-1} \sum_{t=1}^T (e_{ti}^2 - \mathbb{E} e_{ti}^2) \right| &= \max_i \left| T^{-1} \sum_{t=1}^T \sum_{\ell=0}^{L_n} (\varepsilon'_{t-\ell} \phi_{\ell,i} \phi'_{\ell,i} \varepsilon_{t-\ell} - \mathbb{E} \varepsilon'_{t-\ell} \phi_{\ell,i} \phi'_{\ell,i} \varepsilon_{t-\ell}) \right| \\ &\leq T^{-1} \sum_{\ell=0}^{L_n} \max_i |\tilde{\boldsymbol{\varepsilon}}'_\ell \mathbf{A}_{\ell i} \tilde{\boldsymbol{\varepsilon}}_\ell - \mathbb{E} \tilde{\boldsymbol{\varepsilon}}'_\ell \mathbf{A}_{\ell i} \tilde{\boldsymbol{\varepsilon}}_\ell|, \end{aligned}$$

where  $\tilde{\boldsymbol{\varepsilon}}_\ell = (\varepsilon'_{1-\ell}, \dots, \varepsilon'_{T-\ell})' \in \mathbb{R}^{NT}$  and  $\mathbf{A}_{\ell i} = \text{diag}(\phi_{\ell,i} \phi'_{\ell,i}, \dots, \phi_{\ell,i} \phi'_{\ell,i}) \in \mathbb{R}^{NT \times NT}$ . For any  $\ell \in \{0, \dots, L\}$  and  $u > 0$ , the Hanson–Wright inequality in (A.1) with the union bound gives

$$\begin{aligned} \mathbb{P} \left( \max_i |\tilde{\boldsymbol{\varepsilon}}'_\ell \mathbf{A}_{\ell i} \tilde{\boldsymbol{\varepsilon}}_\ell - \mathbb{E} \tilde{\boldsymbol{\varepsilon}}'_\ell \mathbf{A}_{\ell i} \tilde{\boldsymbol{\varepsilon}}_\ell| > u \right) &\leq N \max_i \mathbb{P} (|\tilde{\boldsymbol{\varepsilon}}'_\ell \mathbf{A}_{\ell i} \tilde{\boldsymbol{\varepsilon}}_\ell - \mathbb{E} \tilde{\boldsymbol{\varepsilon}}'_\ell \mathbf{A}_{\ell i} \tilde{\boldsymbol{\varepsilon}}_\ell| > u) \\ &\leq 2N \exp \left( -c \frac{u^2}{K^4 \max_i \|\mathbf{A}_{\ell i}\|_{\text{F}}^2} \right) \end{aligned}$$

Setting  $u = ((\nu+1)/c)^{1/2} K^2 \max_i \|\mathbf{A}_{\ell i}\|_{\text{F}} \log^{1/2}(N \vee T)$  yields

$$\begin{aligned} T^{-1} \sum_{\ell=0}^{L_n} \max_i |\tilde{\boldsymbol{\varepsilon}}'_\ell \mathbf{A}_{\ell i} \tilde{\boldsymbol{\varepsilon}}_\ell - \mathbb{E} \tilde{\boldsymbol{\varepsilon}}'_\ell \mathbf{A}_{\ell i} \tilde{\boldsymbol{\varepsilon}}_\ell| &\leq K^2 T^{-1} \log^{1/2}(N \vee T) \sum_{\ell=0}^{L_n} \max_i \|\mathbf{A}_{\ell i}\|_{\text{F}} \\ &\lesssim T^{-1/2} \log^{1/2}(N \vee T) \sum_{\ell=0}^{L_n} \max_i \|\phi_{\ell,i} \phi'_{\ell,i}\|_{\text{F}} \\ &= T^{-1/2} \log^{1/2}(N \vee T) \sum_{\ell=0}^{\infty} \max_i \|\phi_{\ell,i} \phi'_{\ell,i}\|_2 \\ &\lesssim T^{-1/2} \log^{1/2}(N \vee T) \end{aligned}$$

with probability at least

$$1 - 2N \exp(-(\nu+1) \log(N \vee T)) = 1 - O((N \vee T)^{-\nu}).$$

This completes all the proofs.  $\square$

*Proof of Theorem 1.* We denote by  $\mathbf{M}_{k:\ell} \in \mathbb{R}^{T \times (\ell-k+1)}$  a submatrix of  $\mathbf{M}$  constructed by its  $k$ th to  $\ell$ th columns. Following Ahn and Horenstein (2013), we evaluate the eigenvalues of  $\mathbf{X}\mathbf{X}'$  with recalling notation based on the SVD rather than  $\mathbf{F}^0$  and  $\mathbf{B}^0$ . We define  $\mathbf{P} = \mathbf{V}^0 \mathbf{N}^{-1} \mathbf{V}^{0'}$ ,  $\mathbf{Q} = \mathbf{I}_N - \mathbf{P}$ , and  $\mathbf{U}^* = \mathbf{U}^0 + \mathbf{E} \mathbf{V}^0 \mathbf{N}^{-1} (\mathbf{D}^0)^{-1}$ . Then, we can write  $\mathbf{X}\mathbf{X}' = \mathbf{U}^* \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{*'} + \mathbf{E} \mathbf{Q} \mathbf{E}'$  since  $\mathbf{V}^{0'} \mathbf{V}^0 = \mathbf{N} = \text{diag}(N_1, \dots, N_r)$  by the definition. We also define  $\mathbf{W}_{1:k}$  to be the matrix of  $k$  eigenvectors corresponding to the first  $k$  largest eigenvalues of  $\mathbf{U}^* \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{*'}$ .

We first evaluate the  $r$  largest eigenvalues of  $\mathbf{X}\mathbf{X}'$ . Because  $\lambda_k(\mathbf{U}^0 \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{0'}) = d_k^2 N_k T$ , it is sufficient to show that for any  $k \in \{1, \dots, r\}$ ,

$$\lambda_k(\mathbf{X}\mathbf{X}') = \lambda_k(\mathbf{U}^* \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{*'}) + O(N \vee T), \quad (\text{A.2})$$

$$\lambda_k(\mathbf{U}^* \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{*'}) = \lambda_k(\mathbf{U}^0 \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{0'}) + O\left(T N_1^{1/2} \log^{1/2}(N \vee T) + N \vee T\right). \quad (\text{A.3})$$

Then (A.2) and (A.3) lead to the desired result,

$$\begin{aligned} \lambda_k(\mathbf{X}\mathbf{X}') &= \lambda_k(\mathbf{U}^0 \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{0'}) + O(T N_1^{1/2} \log^{1/2}(N \vee T) + N \vee T) \\ &= d_k^2 N_k T + O\left(T N_1^{1/2} \log^{1/2}(N \vee T) + N \vee T\right). \end{aligned}$$

We show (A.2). Lemma A.5 of Ahn and Horenstein (2013) yields the upper bound

$$\begin{aligned} \sum_{j=1}^k \lambda_j(\mathbf{X}\mathbf{X}') &= \sum_{j=1}^k \lambda_j(\mathbf{U}^* \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{*'} + \mathbf{E} \mathbf{Q} \mathbf{E}') \\ &\leq \sum_{j=1}^k \lambda_j(\mathbf{U}^* \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{*'}) + k \lambda_1(\mathbf{E} \mathbf{Q} \mathbf{E}') \\ &\leq \sum_{j=1}^k \lambda_j(\mathbf{U}^* \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{*'}) + k \lambda_1(\mathbf{E} \mathbf{Q} \mathbf{E}' + \mathbf{E} \mathbf{P} \mathbf{E}') \\ &= \sum_{j=1}^k \lambda_j(\mathbf{U}^* \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{*'}) + k \lambda_1(\mathbf{E} \mathbf{E}'). \end{aligned}$$

Moreover, the lower bound is given by

$$\begin{aligned} \sum_{j=1}^k \lambda_j(\mathbf{X}\mathbf{X}') &\geq T^{-1} \text{tr}(\mathbf{W}_{1:k}' \mathbf{X}\mathbf{X}' \mathbf{W}_{1:k}) \\ &= T^{-1} \text{tr}(\mathbf{W}_{1:k}' \mathbf{U}^* \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{*'} \mathbf{W}_{1:k}) + T^{-1} \text{tr}(\mathbf{W}_{1:k}' \mathbf{E} \mathbf{Q} \mathbf{E}' \mathbf{W}_{1:k}) \\ &\geq \sum_{j=1}^k \lambda_j(\mathbf{U}^* \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{*'}). \end{aligned}$$

Since  $\lambda_1(\mathbf{E} \mathbf{E}') = \|\mathbf{E}\|_2^2 \lesssim T \vee N$  with probability at least  $1 - O((N \vee T)^{-\nu})$  by Lemma 1(a), these two inequalities imply (A.2).

Next, we verify (A.3). By the construction of  $\mathbf{U}^*$ , the upper bound is

$$\begin{aligned} \sum_{j=1}^k \lambda_j(\mathbf{U}^* \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{*'}) &= T^{-1} \text{tr}(\mathbf{W}'_{1:k} \mathbf{U}^0 \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{0'} \mathbf{W}_{1:k}) \\ &\quad + 2T^{-1} \text{tr}(\mathbf{W}'_{1:k} \mathbf{U}^0 \mathbf{D}^0 \mathbf{V}^{0'} \mathbf{E}' \mathbf{W}_{1:k}) + T^{-1} \text{tr}(\mathbf{W}'_{1:k} \mathbf{E} \mathbf{P} \mathbf{E}' \mathbf{W}_{1:k}) \\ &\lesssim \sum_{j=1}^k \lambda_j(\mathbf{U}^0 \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{0'}) + TN_1^{1/2} \log^{1/2}(N \vee T) + N \vee T, \end{aligned}$$

where the last inequality holds by Lemma 3 with probability at least  $1 - O((N \vee T)^{-\nu})$ . Consider the lower bound. We have

$$\begin{aligned} \sum_{j=1}^k \lambda_j(\mathbf{U}^* \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{*'}) &\geq T^{-1} \text{tr}(\mathbf{U}_{1:k}^0 \mathbf{U}^* \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{*'} \mathbf{U}_{1:k}^0) \\ &\gtrsim \sum_{j=1}^k \lambda_j(\mathbf{U}^0 \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{0'}) - TN_1^{1/2} \log^{1/2}(N \vee T). \end{aligned}$$

Finally, we consider the lower and upper bounds of  $\lambda_{r+j}(\mathbf{X} \mathbf{X}')$  for  $j = 1, \dots, k_{\max}$ . Because  $\lambda_{r+j}(\mathbf{U}^* \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{*'}) = 0$  for all  $j \geq 1$ , Lemma 3 entails

$$\begin{aligned} \lambda_{r+j}(\mathbf{X} \mathbf{X}') &\leq \lambda_{r+j}(\mathbf{U}^* \mathbf{D}^0 \mathbf{N} \mathbf{D}^0 \mathbf{U}^{*'}) + \lambda_1(\mathbf{E} \mathbf{Q} \mathbf{E}') \\ &= \lambda_1(\mathbf{E} \mathbf{Q} \mathbf{E}') \lesssim T \vee N \end{aligned}$$

with probability at least  $1 - O((N \vee T)^{-\nu})$ . This completes the proof.  $\square$

*Proof of Theorem 2.* The optimality of the WF-SOFAR estimator implies

$$2^{-1} \|\mathbf{X} - \widehat{\mathbf{F}} \widehat{\mathbf{B}}'\|_{\mathbf{F}}^2 + \eta_n \|\widehat{\mathbf{B}}\|_1 \leq 2^{-1} \|\mathbf{X} - \mathbf{F}^0 \mathbf{B}^{0'}\|_{\mathbf{F}}^2 + \eta_n \|\mathbf{B}^0\|_1.$$

By plugging model (3) and letting  $\Delta = \widehat{\mathbf{F}} \widehat{\mathbf{B}}' - \mathbf{F}^0 \mathbf{B}^{0'}$ , this inequality is equivalently written as

$$2^{-1} \|\mathbf{E} - \Delta\|_{\mathbf{F}}^2 + \eta_n \|\widehat{\mathbf{B}}\|_1 \leq 2^{-1} \|\mathbf{E}\|_{\mathbf{F}}^2 + \eta_n \|\mathbf{B}^0\|_1.$$

Define  $\Delta^f = \widehat{\mathbf{F}} - \mathbf{F}^0$  and  $\Delta^b = \widehat{\mathbf{B}} - \mathbf{B}^0$ . Expanding the first term and using decomposition

$$\Delta = \Delta^f \mathbf{B}^{0'} + \Delta^f \Delta^{b'} + \mathbf{F}^0 \Delta^{b'}$$

lead to

$$\begin{aligned} (1/2) \|\Delta\|_{\mathbf{F}}^2 &\leq \text{tr} \mathbf{E} \Delta' + \eta_n (\|\mathbf{B}^0\|_1 - \|\widehat{\mathbf{B}}\|_1) \\ &\leq \left| \text{tr} \mathbf{E} \mathbf{B}^0 \Delta^{f'} \right| + \left| \text{tr} \mathbf{E} \Delta^b \Delta^{f'} \right| + \left| \text{tr} \Delta^b \mathbf{F}^{0'} \mathbf{E} \right| + \eta_n (\|\mathbf{B}^0\|_1 - \|\widehat{\mathbf{B}}\|_1). \end{aligned} \quad (\text{A.4})$$

We bound the traces in (A.4). By applying Hölder's inequality and using properties of the norms, the first term is bounded as

$$\left| \text{tr} \mathbf{E} \mathbf{B}^0 \Delta^{f'} \right| \leq \|\mathbf{E} \mathbf{B}^0\|_{\max} \|\Delta^f\|_1 \leq (rT)^{1/2} \|\mathbf{E} \mathbf{B}^0\|_{\max} \|\Delta^f\|_{\mathbf{F}}.$$

Similarly, the second and third terms of (A.4) are bounded as

$$\begin{aligned} \left| \text{tr } \mathbf{E} \Delta^b \Delta^{f'} \right| + \left| \text{tr } \Delta^b \mathbf{F}^{0'} \mathbf{E} \right| &\leq \|\mathbf{E} \Delta^b\|_2 \|\Delta^f\|_* + \|\Delta^b\|_1 \|\mathbf{F}^{0'} \mathbf{E}\|_{\max} \\ &\leq r^{1/2} \|\mathbf{E} \Delta^b\|_2 \|\Delta^f\|_F + \|\Delta^b\|_1 \|\mathbf{F}^{0'} \mathbf{E}\|_{\max}. \end{aligned}$$

From these inequalities, the upper bound of (A.4) becomes

$$\begin{aligned} (1/2) \|\Delta\|_F^2 &\leq (rT)^{1/2} \|\mathbf{E} \mathbf{B}^0\|_{\max} \|\Delta^f\|_F + r^{1/2} \|\mathbf{E} \Delta^b\|_2 \|\Delta^f\|_F \\ &\quad + \|\Delta^b\|_1 \|\mathbf{F}^{0'} \mathbf{E}\|_{\max} + \eta_n \left( \|\mathbf{B}^0\|_1 - \|\widehat{\mathbf{B}}\|_1 \right). \end{aligned} \quad (\text{A.5})$$

From Lemmas 1 and 4, there exist some positive constants  $c_1$ – $c_3$  such that the event

$$\begin{aligned} \mathcal{E} = &\left\{ \|\mathbf{E} \Delta^b\|_2 \leq c_1 \|\Delta^b\|_F (\widetilde{N} \vee T)^{1/2} \log^{1/2}(N \vee T) \right\} \\ &\cap \left\{ \|\mathbf{E} \mathbf{B}^0\|_{\max} \leq c_2 N_1^{1/2} \log^{1/2}(N \vee T) \right\} \cap \left\{ \|\mathbf{F}^{0'} \mathbf{E}\|_{\max} \leq c_3 T^{1/2} \log^{1/2}(N \vee T) \right\} \end{aligned}$$

occurs with probability at least  $1 - O((N \vee T)^{-\nu})$  for any fixed constant  $\nu > 0$ . Set the regularization parameter to be  $\eta_n = 2c_3 T^{1/2} \log^{1/2}(N \vee T)$ . Then on event  $\mathcal{E}$ , we have  $\|\mathbf{F}^{0'} \mathbf{E}\|_{\max} \leq \eta_n/2$ , and (A.5) is further bounded as

$$\begin{aligned} \|\Delta\|_F^2 &\lesssim (N_1 T)^{1/2} \log^{1/2}(N \vee T) \|\Delta^f\|_F + (\widetilde{N} \vee T)^{1/2} \log^{1/2}(N \vee T) \|\Delta^b\|_F \|\Delta^f\|_F \\ &\quad + \eta_n \left( \|\Delta^b\|_1 + 2\|\mathbf{B}^0\|_1 - 2\|\widehat{\mathbf{B}}\|_1 \right). \end{aligned} \quad (\text{A.6})$$

We then focus on the last parenthesis of (A.6). Define index set  $\mathcal{S} = \{(i, k) : b_{ik}^0 \neq 0\}$ , the support of  $\mathbf{B}^0$ . Note that  $|\mathcal{S}| = \sum_{k=1}^r N_k \leq rN_1$ . The last parenthesis of (A.6) is rewritten and bounded as

$$\begin{aligned} &\|\Delta^b\|_1 + 2\|\mathbf{B}^0\|_1 - 2\|\widehat{\mathbf{B}}\|_1 \\ &= \|\Delta_{\mathcal{S}}^b\|_1 + \|\Delta_{\mathcal{S}^c}^b\|_1 + 2\|\mathbf{B}_{\mathcal{S}}^0\|_1 - 2\|\widehat{\mathbf{B}}_{\mathcal{S}}\|_1 - 2\|\widehat{\mathbf{B}}_{\mathcal{S}^c}\|_1 \\ &\leq \|\Delta_{\mathcal{S}}^b\|_1 + \|\Delta_{\mathcal{S}^c}^b\|_1 + 2\|\mathbf{B}_{\mathcal{S}}^0\|_1 - 2 \left( \|\mathbf{B}_{\mathcal{S}}^0\|_1 - \|\Delta_{\mathcal{S}}^b\|_1 \right) - 2\|\widehat{\mathbf{B}}_{\mathcal{S}^c}\|_1 \\ &= 3\|\Delta_{\mathcal{S}}^b\|_1 - \|\widehat{\mathbf{B}}_{\mathcal{S}^c}\|_1 \leq 3(rN_1)^{1/2} \|\Delta_{\mathcal{S}}^b\|_F \leq 3(rN_1)^{1/2} \|\Delta^b\|_F. \end{aligned}$$

Therefore, the upper bound of (A.6) is given by

$$\begin{aligned} \|\Delta\|_F^2 &\lesssim (N_1 T)^{1/2} \log^{1/2}(N \vee T) \|\Delta^f\|_F \\ &\quad + (\widetilde{N} \vee T)^{1/2} \log^{1/2}(N \vee T) \|\Delta^b\|_F \|\Delta^f\|_F + N_1^{1/2} \eta_n \|\Delta^b\|_F. \end{aligned} \quad (\text{A.7})$$

Meanwhile, Lemma 5 establishes the lower bound of (A.7). Consequently, we obtain

$$\begin{aligned} \kappa_n \left( \|\Delta^f\|_F^2 + \|\Delta^\lambda\|_F^2 \right) &\lesssim (N_1 T)^{1/2} \log^{1/2}(N \vee T) \|\Delta^f\|_F \\ &\quad + (\widetilde{N} \vee T)^{1/2} \log^{1/2}(N \vee T) \|\Delta^b\|_F \|\Delta^f\|_F + N_1^{1/2} \eta_n \|\Delta^b\|_F \\ &=: \alpha_n \|\Delta^f\|_F + 2\mu_n \|\Delta^b\|_F \|\Delta^f\|_F + \beta_n \|\Delta^b\|_F \\ &\leq \alpha_n \|\Delta^f\|_F + \mu_n \left( \|\Delta^b\|_F^2 + \|\Delta^f\|_F^2 \right) + \beta_n \|\Delta^b\|_F, \end{aligned}$$

where

$$\kappa_n = \frac{N_r(N_r \wedge T)}{N_1},$$

$$\alpha_n = (N_1 T)^{1/2} \log^{1/2}(N \vee T), \quad \mu_n = (\tilde{N} \vee T)^{1/2} \log^{1/2}(N \vee T). \quad \beta_n = N_1^{1/2} \eta_n.$$

Then we have

$$\|\Delta^f\|_F^2 + \|\Delta^b\|_F^2 \leq \frac{(\alpha_n/\kappa_n)\|\Delta^f\|_F + (\beta_n/\kappa_n)\|\Delta^b\|_F}{1 - \mu_n/\kappa_n}.$$

The desired result is obtained by rearranging this inequality. In fact, we have

$$\|\Delta^f\|_F + \|\Delta^b\|_F \leq \frac{3}{2} \left( \frac{\alpha_n/\kappa_n + \beta_n/\kappa_n}{1 - \mu_n/\kappa_n} \right).$$

Finally, since  $\eta_n = 2c_3 T^{1/2} \log^{1/2}(N \vee T)$ , we observe that

$$\begin{aligned} \alpha_n + \beta_n &= (N_1 T)^{1/2} \log^{1/2}(N \vee T) + N_1^{1/2} \eta_n \\ &\lesssim (N_1 T)^{1/2} \log^{1/2}(N \vee T). \end{aligned}$$

This completes the proof.  $\square$

*Proof of Theorem 3.* Following the proof of Theorem 2, we derive the bound. From (A.5) with putting  $\eta_n = 0$ , we have

$$\begin{aligned} (1/2)\|\Delta_{PC}\|_F^2 &\lesssim T^{1/2}\|\mathbf{E}\mathbf{B}^0\|_{\max}\|\Delta_{PC}^f\|_F + \|\mathbf{E}\Delta_{PC}^b\|_2\|\Delta_{PC}^f\|_F + N^{1/2}\|\Delta_{PC}^b\|_F\|\mathbf{F}^{0'}\mathbf{E}\|_{\max}. \end{aligned} \quad (\text{A.8})$$

Lemmas 1 and 4 states that the event

$$\begin{aligned} \mathcal{E} &= \left\{ \|\mathbf{E}\Delta_{PC}^b\|_2 \lesssim \|\Delta_{PC}^b\|_F (N \vee T)^{1/2} \log^{1/2}(N \vee T) \right\} \\ &\cap \left\{ \|\mathbf{E}\mathbf{B}^0\|_{\max} \lesssim N_1^{1/2} \log^{1/2}(N \vee T) \right\} \cap \left\{ \|\mathbf{F}^{0'}\mathbf{E}\|_{\max} \lesssim T^{1/2} \log^{1/2}(N \vee T) \right\} \end{aligned}$$

occurs with probability at least  $1 - O((N \vee T)^{-\nu})$  for any fixed constant  $\nu > 0$ . On event  $\mathcal{E}$  together with Lemma 5, (A.8) becomes

$$\kappa_n \left( \|\Delta_{PC}^f\|_F^2 + \|\Delta_{PC}^b\|_F^2 \right) \leq \alpha_n \|\Delta_{PC}^f\|_F + \mu_n \left( \|\Delta_{PC}^b\|_F^2 + \|\Delta_{PC}^f\|_F^2 \right) + \beta_n \|\Delta_{PC}^b\|_F,$$

where

$$\begin{aligned} \kappa_n &= \frac{N_r(N_r \wedge T)}{N_1}, \quad \mu_n = (N \vee T)^{1/2} \log^{1/2}(N \vee T) \\ \alpha_n &= (N_1 T)^{1/2} \log^{1/2}(N \vee T), \quad \beta_n = (NT)^{1/2} \log^{1/2}(N \vee T). \end{aligned}$$

The desired result is obtained by rearranging this inequality as in the proof of Theorem 2. In fact, we have

$$\|\Delta_{PC}^f\|_F + \|\Delta_{PC}^b\|_F \leq \frac{3}{2} \left( \frac{\alpha_n/\kappa_n + \beta_n/\kappa_n}{1 - \mu_n/\kappa_n} \right).$$

Finally, we observe that

$$\begin{aligned}\alpha_n + \beta_n &= (N_1 T)^{1/2} \log^{1/2}(N \vee T) + (NT)^{1/2} \log^{1/2}(N \vee T) \\ &\lesssim (NT)^{1/2} \log^{1/2}(N \vee T).\end{aligned}$$

This completes the proof of Theorem 3.  $\square$

*Proof of Theorem 4.* Throughout this proof, we omit the superscript of the adaptive estimators  $(\widehat{\mathbf{F}}^{\text{ada}}, \widehat{\mathbf{B}}^{\text{ada}})$  and simply write them as  $(\widehat{\mathbf{F}}, \widehat{\mathbf{B}})$ . Recall  $\mathcal{S} = \text{supp}(\mathbf{B}^0)$ , which is a subset of  $\{1, \dots, N\} \times \{1, \dots, r\}$ . For any matrix  $\mathbf{B} = (b_{ik}) \in \mathbb{R}^{N \times r}$ , define  $\mathbf{B}_{\mathcal{S}} \in \mathbb{R}^{N \times r}$  as the matrix whose  $(i, k)$ th element is  $b_{ik} 1\{(i, k) \in \mathcal{S}\}$ . Similarly, define  $\mathbf{B}_{\mathcal{S}^c} \in \mathbb{R}^{N \times r}$  whose  $(i, k)$ th element is  $b_{ik} 1\{(i, k) \in \mathcal{S}^c\}$ . By the definition, note that  $\mathbf{B}_{\mathcal{S}}^0 = \mathbf{B}^0$  and  $\mathbf{B}_{\mathcal{S}^c}^0 = \mathbf{0}$ . Recall that the objective function for obtaining the adaptive WF-SOFAR estimator is given by

$$Q_n(\mathbf{F}, \mathbf{B}) := \frac{1}{2} \|\mathbf{X} - \mathbf{F}\mathbf{B}'\|_{\mathbf{F}}^2 + \eta_n \|\mathbf{W} \circ \mathbf{B}\|_1 \quad (\text{A.9})$$

subject to  $\mathbf{F}'\mathbf{F}/T = \mathbf{I}_r$  and  $\mathbf{B}'\mathbf{B}$  being diagonal. The strategy of this proof consists of two steps. In the first step, we show that the *oracle estimator*  $(\widehat{\mathbf{F}}^o, \widehat{\mathbf{B}}_{\mathcal{S}}^o)$ , which is defined as a minimizer of  $Q_n(\mathbf{F}, \mathbf{B}_{\mathcal{S}})$  (i.e., a minimizer of the correctly zero-restricted minimization problem), is consistent to  $(\mathbf{F}^0, \mathbf{B}_{\mathcal{S}}^0)$  with some rate of convergence. In the second step, we prove that the oracle estimator is indeed a minimizer of the unrestricted problem,  $\min Q_n(\mathbf{F}, \mathbf{B})$  over  $\mathbb{R}^{T \times r} \times \mathbb{R}^{N \times r}$ .

(First step) We derive the rate of convergence of the oracle estimator. To this end, it suffices to show that as  $n \rightarrow \infty$ , there exists a (large) constant  $C > 0$  such that

$$\mathbb{P} \left( \inf_{\|\mathbf{U}\|_{\mathbf{F}}=C, \|\mathbf{V}_{\mathcal{S}}\|_{\mathbf{F}}=C} Q_n(\mathbf{F}^0 + r_n \mathbf{U}, \mathbf{B}_{\mathcal{S}}^0 + r_n \mathbf{V}_{\mathcal{S}}) > Q_n(\mathbf{F}^0, \mathbf{B}_{\mathcal{S}}^0) \right) \rightarrow 1, \quad (\text{A.10})$$

where  $\mathbf{U} \in \mathbb{R}^{T \times r}$  and  $\mathbf{V} \in \mathbb{R}^{N \times r}$  are deterministic matrices, and

$$r_n = \frac{N_1(N_1 T)^{1/2} \log^{1/2}(N \vee T)}{N_r(N_r \wedge T)} = \frac{\gamma_n(\tilde{N}) N_1^{1/2} T^{1/2}}{(\tilde{N} \vee T)^{1/2}}.$$

This implies that the oracle estimator  $(\widehat{\mathbf{F}}^o, \widehat{\mathbf{B}}_{\mathcal{S}}^o)$  lies in the ball

$$\{(\mathbf{F}, \mathbf{B}_{\mathcal{S}}) \in \mathbb{R}^{T \times r} \times \mathbb{R}^{N \times r} : \|\mathbf{F} - \mathbf{F}^0\|_{\mathbf{F}} \leq C r_n, \|\mathbf{B}_{\mathcal{S}} - \mathbf{B}_{\mathcal{S}}^0\|_{\mathbf{F}} \leq C r_n\}$$

with high probability, which gives the desired rate of convergence.

To show (A.10), we have

$$\begin{aligned}& Q_n(\mathbf{F}^0 + r_n \mathbf{U}, \mathbf{B}_{\mathcal{S}}^0 + r_n \mathbf{V}_{\mathcal{S}}) - Q_n(\mathbf{F}^0, \mathbf{B}_{\mathcal{S}}^0) \\ &= 2^{-1} \|\mathbf{X} - (\mathbf{F}^0 + r_n \mathbf{U})(\mathbf{B}_{\mathcal{S}}^0 + r_n \mathbf{V}_{\mathcal{S}})'\|_{\mathbf{F}}^2 - 2^{-1} \|\mathbf{X} - \mathbf{F}^0 \mathbf{B}_{\mathcal{S}}^0\|_{\mathbf{F}}^2 \\ &\quad + \eta_n \|\mathbf{W} \circ (\mathbf{B}_{\mathcal{S}}^0 + r_n \mathbf{V}_{\mathcal{S}})\|_1 - \eta_n \|\mathbf{W} \circ \mathbf{B}_{\mathcal{S}}^0\|_1 \\ &\geq \text{tr } \mathbf{E}'(r_n \mathbf{F}^0 \mathbf{V}_{\mathcal{S}}' + r_n \mathbf{U} \mathbf{B}_{\mathcal{S}}^0{}' + r_n^2 \mathbf{U} \mathbf{V}_{\mathcal{S}}') \\ &\quad + 2^{-1} \|r_n \mathbf{F}^0 \mathbf{V}_{\mathcal{S}}' + r_n \mathbf{U} \mathbf{B}_{\mathcal{S}}^0{}' + r_n^2 \mathbf{U} \mathbf{V}_{\mathcal{S}}'\|_{\mathbf{F}}^2 - r_n \eta_n \|\mathbf{W}_{\mathcal{S}} \circ \mathbf{V}_{\mathcal{S}}\|_1 \\ &= (I) + (II) + (III).\end{aligned} \quad (\text{A.11})$$



By Lemma 7 (a)–(c), we bound  $(I)$  as

$$\begin{aligned}
|(I)| &= \left| r_n \operatorname{tr} \mathbf{E}' \mathbf{F}^0 \mathbf{V}'_{\mathcal{S}} + r_n \operatorname{tr} \mathbf{E}' \mathbf{U} \mathbf{B}_{\mathcal{S}}^{0'} + r_n^2 \operatorname{tr} \mathbf{E}' \mathbf{U} \mathbf{V}'_{\mathcal{S}} \right| \\
&\leq r_n \left| \operatorname{tr} \mathbf{V}'_{\mathcal{S}} \mathbf{E}' \mathbf{F}^0 \right| + r_n \left| \operatorname{tr} \mathbf{B}_{\mathcal{S}}^{0'} \mathbf{E}' \mathbf{U} \right| + r_n^2 \left| \operatorname{tr} \mathbf{V}'_{\mathcal{S}} \mathbf{E}' \mathbf{U} \right| \\
&\lesssim r_n \left( T^{1/2} \|\mathbf{V}_{\mathcal{S}}\|_{\mathbb{F}} + N_1^{1/2} \|\mathbf{U}\|_{\mathbb{F}} \right) \log^{1/2}(N \vee T) + r_n^2 \|\mathbf{U}\|_{\mathbb{F}} \|\mathbf{V}_{\mathcal{S}}\|_{\mathbb{F}} \log^{1/2}(N \vee T).
\end{aligned}$$

Next, we bound  $(II)$  from below as

$$\begin{aligned}
(II) &= 2^{-1} \|\mathbf{U} \mathbf{B}^{0'} + \mathbf{U} \mathbf{V}' + \mathbf{F}^0 \mathbf{V}'\|_{\mathbb{F}}^2 \\
&= 2^{-1} \|\mathbf{U} \mathbf{B}^{0'}\|_{\mathbb{F}}^2 + 2^{-1} \|\mathbf{U} \mathbf{V}'\|_{\mathbb{F}}^2 + 2^{-1} \|\mathbf{F}^0 \mathbf{V}'\|_{\mathbb{F}}^2 \\
&\quad + \operatorname{tr} \mathbf{V} \mathbf{U}' \mathbf{F}^0 \mathbf{V}' + \operatorname{tr} \mathbf{B}^0 \mathbf{U}' \mathbf{U} \mathbf{V}' + \operatorname{tr} \mathbf{B}^0 \mathbf{U}' \mathbf{F}^0 \mathbf{V}' \\
&\geq 2^{-1} \|\mathbf{U} \mathbf{B}^{0'}\|_{\mathbb{F}}^2 + 2^{-1} \|\mathbf{F}^0 \mathbf{V}'\|_{\mathbb{F}}^2 - |\operatorname{tr} \mathbf{V} \mathbf{U}' \mathbf{F}^0 \mathbf{V}'| - |\operatorname{tr} \mathbf{B}^0 \mathbf{U}' \mathbf{U} \mathbf{V}'| - |\operatorname{tr} \mathbf{B}^0 \mathbf{U}' \mathbf{F}^0 \mathbf{V}'| \\
&= (i) + (ii) + (iii) + (iv) + (v).
\end{aligned}$$

In view of the Rayleigh quotient,  $(i)$  and  $(ii)$  are further bounded from below as

$$\begin{aligned}
(i) + (ii) &= 2^{-1} \|\mathbf{U} \mathbf{B}^{0'}\|_{\mathbb{F}}^2 + 2^{-1} \|\mathbf{F}^0 \mathbf{V}'_{\mathcal{S}}\|_{\mathbb{F}}^2 \\
&= 2^{-1} r_n^2 \|(\mathbf{I}_T \otimes \mathbf{B}^0) \operatorname{vec}(\mathbf{U}')\|_2^2 + 2^{-1} r_n^2 \|(\mathbf{I}_N \otimes \mathbf{F}^0) \operatorname{vec}(\mathbf{V}'_{\mathcal{S}})\|_2^2 \\
&\gtrsim r_n^2 \left\{ \min_{\mathbf{u} \in \mathbb{R}^{T_r} \setminus \{\mathbf{0}\}} \left( \frac{\|(\mathbf{I}_T \otimes \mathbf{B}^0) \mathbf{u}\|_2^2}{\|\mathbf{u}\|_2^2} \right) \|\mathbf{U}\|_{\mathbb{F}}^2 + \min_{\mathbf{v} \in \mathbb{R}^{N_r} \setminus \{\mathbf{0}\}} \left( \frac{\|(\mathbf{I}_N \otimes \mathbf{F}^0) \mathbf{v}\|_2^2}{\|\mathbf{v}\|_2^2} \right) \|\mathbf{V}_{\mathcal{S}}\|_{\mathbb{F}}^2 \right\} \\
&\gtrsim r_n^2 (N_r \|\mathbf{U}\|_{\mathbb{F}}^2 + T \|\mathbf{V}_{\mathcal{S}}\|_{\mathbb{F}}^2).
\end{aligned}$$

Meanwhile, by Lemma 7 (d)–(f),  $|(iii) + (iv) + (v)|$  is bounded from above as

$$\begin{aligned}
|(iii) + (iv) + (v)| &\lesssim r_n^3 \left( \|\mathbf{U}\|_{\mathbb{F}} \|\mathbf{V}_{\mathcal{S}}\|_{\mathbb{F}}^2 \log^{1/2}(N \vee T) + N_1^{1/2} \|\mathbf{U}\|_{\mathbb{F}}^2 \|\mathbf{V}_{\mathcal{S}}\|_{\mathbb{F}} \right) \\
&\quad + r_n^2 N_1^{1/2} \|\mathbf{U}\|_{\mathbb{F}} \|\mathbf{V}_{\mathcal{S}}\|_{\mathbb{F}} \log^{1/2}(N \vee T).
\end{aligned}$$

Combining these bounds of  $(i)$ – $(v)$  yields

$$\begin{aligned}
(II) &\gtrsim (i) + (ii) - |(iii) + (iv) + (v)| \\
&\gtrsim r_n^2 (N_r \|\mathbf{U}\|_{\mathbb{F}}^2 + T \|\mathbf{V}_{\mathcal{S}}\|_{\mathbb{F}}^2) \\
&\quad - r_n^3 \left( \|\mathbf{U}\|_{\mathbb{F}} \|\mathbf{V}_{\mathcal{S}}\|_{\mathbb{F}}^2 \log^{1/2}(N \vee T) + N_1^{1/2} \|\mathbf{U}\|_{\mathbb{F}}^2 \|\mathbf{V}_{\mathcal{S}}\|_{\mathbb{F}} \right) \\
&\quad - r_n^2 N_1^{1/2} \|\mathbf{U}\|_{\mathbb{F}} \|\mathbf{V}_{\mathcal{S}}\|_{\mathbb{F}} \log^{1/2}(N \vee T).
\end{aligned}$$

We then consider  $(III)$  in (A.11). Lemma 8 yields

$$|(III)| = r_n \eta_n \|\mathbf{W}_{\mathcal{S}} \circ \mathbf{V}_{\mathcal{S}}\|_1 \leq r_n \eta_n \|\mathbf{W}_{\mathcal{S}}\|_{\mathbb{F}} \|\mathbf{V}_{\mathcal{S}}\|_{\mathbb{F}} \lesssim N_1^{1/2} r_n (\eta_n / \underline{b}_n^0) \|\mathbf{V}_{\mathcal{S}}\|_{\mathbb{F}}.$$

Putting together the pieces obtained so far with (A.11), we have

$$\begin{aligned}
& \inf_{\|\mathbf{U}\|_F=C, \|\mathbf{V}_S\|_F=C} Q_n(\mathbf{F}^0 + r_n \mathbf{U}, \mathbf{B}_S^0 + r_n \mathbf{V}_S) - Q_n(\mathbf{F}^0, \mathbf{B}_S^0) \\
& \gtrsim \inf_{\|\mathbf{U}\|_F=C, \|\mathbf{V}_S\|_F=C} \{(II) - |(I)| - |(III)|\} \\
& \gtrsim \inf_{\|\mathbf{U}\|_F=C, \|\mathbf{V}_S\|_F=C} \left\{ r_n^2 (N_r \|\mathbf{U}\|_F^2 + T \|\mathbf{V}_S\|_F^2) \right. \\
& \quad - r_n^3 \left( \|\mathbf{U}\|_F \|\mathbf{V}_S\|_F^2 \log^{1/2}(N \vee T) + N_1^{1/2} \|\mathbf{U}\|_F^2 \|\mathbf{V}_S\|_F \right) \\
& \quad - r_n^2 N_1^{1/2} \|\mathbf{U}\|_F \|\mathbf{V}_S\|_F \log^{1/2}(N \vee T) \\
& \quad - r_n \left( T^{1/2} \|\mathbf{V}_S\|_F \log^{1/2}(N \vee T) + N_1^{1/2} \|\mathbf{U}\|_F \log^{1/2}(N \vee T) \right) \\
& \quad \left. - r_n^2 \|\mathbf{U}\|_F \|\mathbf{V}_S\|_F \log^{1/2}(N \vee T) - N_1^{1/2} r_n (\eta_n / \underline{b}_n^0) \|\mathbf{V}_S\|_F \right\} \\
& \sim r_n^2 \left( N_r + T - N_1^{1/2} \log^{1/2}(N \vee T) \right) C^2 - r_n^3 N_1^{1/2} C^3 \\
& \quad - r_n \left( T^{1/2} \log^{1/2}(N \vee T) + N_1^{1/2} \log^{1/2}(N \vee T) + N_1^{1/2} (\eta_n / \underline{b}_n^0) \right) C.
\end{aligned}$$

By condition (8) and the fact that

$$r_n = \frac{N_1}{N_r} \cdot \frac{N_1^{1/2} T^{1/2} \log^{1/2}(N \vee T)}{N_r \wedge T} \geq 1 \cdot \frac{N_r^{1/2} T^{1/2}}{N_r \wedge T} \geq 1,$$

we have

$$\begin{aligned}
& \inf_{\|\mathbf{U}\|_F=C, \|\mathbf{V}_S\|_F=C} Q_n(\mathbf{F}^0 + r_n \mathbf{U}, \mathbf{B}_S^0 + r_n \mathbf{V}_S) - Q_n(\mathbf{F}^0, \mathbf{B}_S^0) \\
& \gtrsim r_n^2 (N_r \vee T) C^2 - r_n^3 N_1^{1/2} C^3 - r_n N_1^{1/2} (\eta_n / \underline{b}_n^0) C.
\end{aligned} \tag{A.12}$$

Furthermore, in (A.12), the first term asymptotically dominates the second term because the ratio,

$$\frac{r_n^3 N_1^{1/2}}{r_n^2 (N_r \vee T)} = \frac{N_1^2}{N_r^2 T^{1/2}},$$

converges to zero by condition (15). Finally, we compare the first term with the third term in (A.12). Observe that conditions (17) (first bound) and (18) give

$$\underline{b}_n^0 \gtrsim \frac{N_1^{1/2} \eta_n}{r_n (N_r \vee T)}.$$

Therefore, the first term in (A.12) becomes strictly larger than the third term in absolute value as long as  $C > 0$  is taken to be large enough. This means that the lower bound tends to positive for such  $C > 0$  and (A.10) holds.

(Second step) Set  $\widehat{\mathbf{F}} = \widehat{\mathbf{F}}^o$  and  $\widehat{\mathbf{B}} = \widehat{\mathbf{B}}_S^o$ . If the estimator  $(\widehat{\mathbf{F}}, \widehat{\mathbf{B}})$  is indeed a minimizer of the unrestricted problem,  $\min Q_n(\mathbf{F}, \mathbf{B})$  over  $\mathbb{R}^{T \times r} \times \mathbb{R}^{N \times r}$ , the proof completes. Note that  $\text{supp } \widehat{\mathbf{B}} = \mathcal{S}$  by the construction. taking the same strategy as in Fan et al. (2014), we check the optimality of  $(\widehat{\mathbf{F}}, \widehat{\mathbf{B}})$ . By a simple calculation, the (sub-)gradients of  $Q_n$  with respect to  $\mathbf{F}$  and  $\mathbf{B}$  are given by

$$\begin{aligned}
\nabla_{\mathbf{F}} Q_n(\mathbf{F}, \mathbf{B}) &= \mathbf{F} \mathbf{B}' \mathbf{B} - \mathbf{X} \mathbf{B}, \\
\nabla_{\mathbf{B}} Q_n(\mathbf{F}, \mathbf{B}) &= \mathbf{B} \mathbf{F}' \mathbf{F} - \mathbf{X}' \mathbf{F} + \eta_n \mathbf{T},
\end{aligned}$$

where the  $(i, k)$ th element of  $\mathbf{T} \in \mathbb{R}^{N \times r}$  is defined as

$$t_{ik} \begin{cases} = w_{ik} \operatorname{sgn}(b_{ik}) & \text{for } b_{ik} \neq 0, \\ \in w_{ik}[-1, 1] & \text{for } b_{ik} = 0. \end{cases}$$

Then  $(\widehat{\mathbf{F}}, \widehat{\mathbf{B}})$  is a strict minimizer of (6) if the following conditions hold:

$$\widehat{\mathbf{F}}\widehat{\mathbf{B}}'\widehat{\mathbf{B}} - \mathbf{X}\widehat{\mathbf{B}} = \mathbf{0}_{T \times r}, \quad (\text{A.13})$$

$$T\widehat{\mathbf{B}}_{\mathcal{S}} - (\mathbf{X}'\widehat{\mathbf{F}})_{\mathcal{S}} + \eta_n \mathbf{W}_{\mathcal{S}} \circ \operatorname{sgn} \widehat{\mathbf{B}}_{\mathcal{S}} = \mathbf{0}_{N \times r}, \quad (\text{A.14})$$

$$\left\| \mathbf{W}_{\mathcal{S}^c}^- \circ \left\{ T\widehat{\mathbf{B}}_{\mathcal{S}^c} - (\mathbf{X}'\widehat{\mathbf{F}})_{\mathcal{S}^c} \right\} \right\|_{\max} < \eta_n, \quad (\text{A.15})$$

where  $\widehat{\mathbf{F}}'\widehat{\mathbf{F}} = T\mathbf{I}_r$  has been used, and  $\mathbf{W}^- \in \mathbb{R}^{N \times r}$  is the matrix with its  $(i, k)$ th elements given by  $1/w_{ik}$ . Since  $(\widehat{\mathbf{F}}, \widehat{\mathbf{B}}_{\mathcal{S}})$  is a minimizer of  $Q_n(\mathbf{F}, \mathbf{B}_{\mathcal{S}})$ , it satisfies the Karush–Kuhn–Tucker (KKT) conditions. Therefore, we only need to check condition (A.15). Condition (A.15) is indeed verified by Lemma 9. This completes the proof of Theorem 4.  $\square$

*Proof of Corollary 4.* Recall that  $\hat{\alpha}_j = \log \widehat{N}_j / \log N$  with  $\widehat{N}_j = |\operatorname{supp}(\widehat{\mathbf{b}}_j^{\text{ada}})|$  and  $\alpha_j = \log N_j / \log N$  by the definition. Because  $\{\operatorname{supp}(\widehat{\mathbf{B}}^{\text{ada}}) = \operatorname{supp}(\mathbf{B}^0)\} \subset \{\widehat{N}_j = N_j \text{ for all } j = 1, \dots, r\}$ , we have

$$\begin{aligned} \mathbb{P}(\hat{\alpha}_j = \alpha_j \text{ for all } j = 1, \dots, r) &= \mathbb{P}(\widehat{N}_j = N_j \text{ for all } j = 1, \dots, r) \\ &\geq \mathbb{P}(\operatorname{supp}(\widehat{\mathbf{B}}^{\text{ada}}) = \operatorname{supp}(\mathbf{B}^0)). \end{aligned}$$

The last probability tends to one by the factor selection consistency. This completes the proof of Corollary 4.  $\square$

## B Related Lemmas and their Proofs

**Lemma 2.** Assume  $X_i \sim \text{ind. subG}(\alpha_i^2)$  and  $Y_i \sim \text{ind. subE}(\gamma_i)$ . Then, for any deterministic sequences  $(\phi_i)$  and  $(\psi_i)$ , the following statements are true:

- (a)  $X_i X_j \sim \text{subE}(4e\alpha_i \alpha_j)$  for  $i \neq j$ .
- (b)  $\sum_{i=1}^n \phi_i X_i \sim \text{subG}(\sum_{i=1}^n \phi_i^2 \alpha_i^2)$ .
- (c)  $\sum_{i=1}^n \psi_i Y_i \sim \text{subE}((\sum_{i=1}^n \psi_i^2 \gamma_i^2)^{1/2}, \max_i |\psi_i| \gamma_i)$ .

*Proof.* This proof was achieved in Uematsu and Tanaka (2019), but is repeated here for completeness. (a) Since  $X_i$  is  $\text{subG}(\alpha_i^2)$ , we obtain  $\mathbb{E}|X|^k \leq (2\alpha_i^2)^{k/2} k\Gamma(k/2)$ ; see Rigollet and Hütter (2017), for instance. Then we see from the dominated convergence theorem and

independence that

$$\begin{aligned}
\mathbb{E} \exp(sX_i X_j) &= 1 + \sum_{k=2}^{\infty} \frac{s^k \mathbb{E}(X_i X_j)^k}{k!} \leq 1 + \sum_{k=2}^{\infty} \frac{s^k \mathbb{E}|X_i|^k \mathbb{E}|X_j|^k}{k!} \\
&\leq 1 + \sum_{k=2}^{\infty} \frac{s^k (2\alpha_i \alpha_j)^k k^2 \Gamma(k/2)^2}{k!} \leq 1 + \sum_{k=2}^{\infty} \frac{s^k (2\alpha_i \alpha_j)^k k^2 (k/2)^k}{k!} \\
&= 1 + \sum_{k=2}^{\infty} \frac{s^k (\alpha_i \alpha_j)^k k^{k+2}}{k!} \leq 1 + \sum_{k=2}^{\infty} (2e\alpha_i \alpha_j s)^k \\
&= 1 + (2e\alpha_i \alpha_j s)^2 \sum_{k=0}^{\infty} (2e\alpha_i \alpha_j s)^k,
\end{aligned}$$

where we have used  $\Gamma(k/2) \leq (k/2)^{k/2}$  and  $k^{k+2} \leq (2\pi)^{-1/2} k! e^k k^{3/2} \leq k! (2e)^k$ . Therefore, for any  $|s| \leq (4e\alpha_i \alpha_j)^{-1}$ , it holds that

$$\mathbb{E} \exp(sX_i X_j) \leq 1 + 8(e\alpha_i \alpha_j)^2 s^2 \leq \exp((4e\alpha_i \alpha_j)^2 s^2 / 2).$$

This means that the product  $X_i X_j$  is subE( $4e\alpha_i \alpha_j$ ).

(b) By the definition of subG, we have

$$\begin{aligned}
\mathbb{E} \exp\left(s \sum_{i=1}^n \phi_i X_i\right) &= \prod_{i=1}^n \mathbb{E} \exp(s \phi_i X_i) \\
&\leq \prod_{i=1}^n \exp(s^2 \phi_i^2 \alpha_i^2 / 2) = \exp\left(s^2 \sum_{i=1}^n \phi_i^2 \alpha_i^2 / 2\right),
\end{aligned}$$

which yields the result.

(c) First note that  $\psi_i Y_i \sim \text{subE}(\psi_i \gamma_i, |\psi_i| \gamma_i)$  because  $\mathbb{E} \exp(s \psi_i Y_i) \leq \exp(s^2 \psi_i^2 \gamma_i^2 / 2)$  holds for all  $|s| \leq (|\psi_i| \gamma_i)^{-1}$ . Thus, we can see that

$$\begin{aligned}
\mathbb{E} \exp\left(s \sum_{i=1}^n \psi_i Y_i\right) &= \prod_{i=1}^n \mathbb{E} \exp(s \psi_i Y_i) \\
&\leq \prod_{i=1}^n \exp(s^2 \psi_i^2 \gamma_i^2 / 2) = \exp\left(s^2 \sum_{i=1}^n \psi_i^2 \gamma_i^2 / 2\right),
\end{aligned}$$

where the inequality holds for all  $|s| \leq (\max_i |\psi_i| \gamma_i)^{-1}$ . This gives the result by the definition of subE, and completes all the proofs.  $\square$

**Lemma 3.** Suppose the same conditions as Theorem 1. Then, for any  $\mathbf{H} \in \mathbb{R}^{T \times k}$  ( $k \leq r$ ) such that  $\mathbf{H}'\mathbf{H} = T\mathbf{I}_k$ , the following inequalities simultaneously hold with probability at least  $1 - O((N \vee T)^{-\nu})$ :

- (a)  $T^{-1} \left| \text{tr} \mathbf{H}' \mathbf{U}^0 \mathbf{D}^0 \mathbf{V}^{0'} \mathbf{E}' \mathbf{H} \right| \lesssim T N_1^{1/2} \log^{1/2}(N \vee T),$
- (b)  $T^{-1} \text{tr} \mathbf{H}' \mathbf{E} \mathbf{P} \mathbf{E}' \mathbf{H} \lesssim N \vee T,$
- (c)  $\lambda_1(\mathbf{E} \mathbf{Q} \mathbf{E}') \lesssim T \vee N,$
- (d)  $T^{-1} \text{tr}(\mathbf{H}' \mathbf{E} \mathbf{Q} \mathbf{E}' \mathbf{H}) \lesssim T \vee N.$

*Proof.* Recall the notation based on the SVD of  $\mathbf{C}^0$ :  $\mathbf{U}^0 = \mathbf{F}^0$  and  $\mathbf{V}^0 \mathbf{D}^0 = \mathbf{B}^0$ . We derive the results on the event that Lemma 1 hold, which occurs with probability at least  $1 - O((N \vee T)^{-\nu})$ . Prove (a). Low rankness of each matrix and Lemma 1(b) give

$$\begin{aligned} \left| \text{tr} \mathbf{H}' \mathbf{U}^0 \mathbf{D}^0 \mathbf{V}^{0'} \mathbf{E}' \mathbf{H} \right| &\leq \|\mathbf{H} \mathbf{H}'\|_F \|\mathbf{U}^0\|_F \|\mathbf{D}^0 \mathbf{V}^{0'} \mathbf{E}'\|_F \\ &\lesssim \|\mathbf{H} \mathbf{H}'\|_F \|\mathbf{U}^0\|_F \|\mathbf{D}^0 \mathbf{V}^{0'} \mathbf{E}'\|_2 \\ &\lesssim T T^{1/2} T^{1/2} \|\mathbf{D}^0 \mathbf{V}^{0'} \mathbf{E}'\|_{\max} \\ &\lesssim T^2 N_1^{1/2} \log^{1/2}(N \vee T). \end{aligned}$$

Prove (b). Since the rank of  $\mathbf{P}$  is at most  $r$ , Lemma 1(a) gives

$$\text{tr} \mathbf{H}' \mathbf{E} \mathbf{P} \mathbf{E}' \mathbf{H} \lesssim \|\mathbf{H} \mathbf{H}'\|_F \|\mathbf{E} \mathbf{P} \mathbf{E}'\|_2 \leq T \|\mathbf{E}\|_2^2 \|\mathbf{P}\|_2 \lesssim T(N \vee T).$$

Prove (c). By the argument of the proof of Lemma A.8 in Ahn and Horenstein (2013) and Lemma 1(a), the bound

$$\lambda_1(\mathbf{E} \mathbf{Q} \mathbf{E}') \leq \lambda_1(\mathbf{E} \mathbf{Q} \mathbf{E}' + \mathbf{E} \mathbf{P} \mathbf{E}') = \lambda_1(\mathbf{E} \mathbf{E}') = \|\mathbf{E}\|_2^2 \lesssim T \vee N.$$

Prove (d). From the triangle inequality and result (c), we have

$$\text{tr}(\mathbf{H}' \mathbf{E} \mathbf{Q} \mathbf{E}' \mathbf{H}) \lesssim \|\mathbf{H} \mathbf{H}'\|_F \|\mathbf{E} \mathbf{Q} \mathbf{E}'\|_2 \leq \|\mathbf{H} \mathbf{H}'\|_F (\|\mathbf{E} \mathbf{E}'\|_2 + \|\mathbf{E} \mathbf{P} \mathbf{E}'\|_2) \lesssim T(T \vee N).$$

This completes all the proofs of (a)–(d).  $\square$

**Lemma 4.** *Suppose the same conditions as Theorem 2. Then we have*

$$\|\mathbf{E} \mathbf{\Delta}^b\|_2 \lesssim \|\mathbf{\Delta}^b\|_F (\tilde{N} \vee T)^{1/2} \log^{1/2}(N \vee T)$$

*with probability at least  $1 - O((N \vee T)^{-\nu})$ .*

*Proof.* In the upper bound of (A.4), we consider a tighter bound of the second trace. The second trace in the upper bound of (A.4) is bounded as

$$\left| \text{tr} \mathbf{E} \mathbf{\Delta}^b \mathbf{\Delta}^{f'} \right| \leq \|\mathbf{E} \mathbf{\Delta}^b\|_2 \|\mathbf{\Delta}^f\|_*.$$

Because  $\hat{\mathbf{B}}$  and  $\mathbf{B}^0$  lie in the set  $\mathcal{B}(\tilde{N}) = \{\mathbf{B} \in \mathbb{R}^{N \times r} : \|\mathbf{B}\|_0 \lesssim \tilde{N}/2\}$  for  $\tilde{N} \in [N_1, N]$  by Assumption 4, we have

$$\|\mathbf{\Delta}^b\|_0 \leq \|\hat{\mathbf{B}}\|_0 + \|\mathbf{B}^0\|_0 \lesssim \tilde{N}/2 + \tilde{N}/2 \leq \tilde{N}.$$

Define a set of sparse vectors  $\mathcal{V}(\mathcal{A}) = \{\mathbf{v} \in \mathbb{R}^N \setminus \{\mathbf{0}\} : \|\mathbf{v}\|_0 = |\mathcal{A}|\}$  with  $\mathcal{A} \subset \{1, \dots, N\}$ . Then, by the definition of the spectral norm, we have

$$\begin{aligned} \|\mathbf{E} \mathbf{\Delta}^b\|_2^2 &= \max_{\mathbf{u} \in \mathbb{R}^r \setminus \{\mathbf{0}\}} \frac{\mathbf{u}' \mathbf{\Delta}^b \mathbf{E}' \mathbf{E} \mathbf{\Delta}^b \mathbf{u}}{\mathbf{u}' \mathbf{u}} \\ &\leq \max_{\mathbf{u} \in \mathbb{R}^r \setminus \{\mathbf{0}\}} \frac{\mathbf{u}' \mathbf{\Delta}^b \mathbf{E}' \mathbf{E} \mathbf{\Delta}^b \mathbf{u}}{\mathbf{u}' \mathbf{\Delta}^b \mathbf{E}' \mathbf{E} \mathbf{\Delta}^b \mathbf{u}} \max_{\mathbf{u} \in \mathbb{R}^r \setminus \{\mathbf{0}\}} \frac{\mathbf{u}' \mathbf{\Delta}^b \mathbf{\Delta}^b \mathbf{u}}{\mathbf{u}' \mathbf{u}} \\ &\leq \max_{|\mathcal{A}| \lesssim \tilde{N}} \max_{\mathbf{v} \in \mathcal{V}(\mathcal{A})} \frac{\mathbf{v}' \mathbf{E}' \mathbf{E} \mathbf{v}}{\mathbf{v}' \mathbf{v}} \|\mathbf{\Delta}^b\|_2^2 = \max_{|\mathcal{A}| \lesssim \tilde{N}} \max_{\mathbf{v}_{\mathcal{A}} \in \mathbb{R}^{|\mathcal{A}|}} \frac{\mathbf{v}'_{\mathcal{A}} \mathbf{E}'_{\mathcal{A}} \mathbf{E}_{\mathcal{A}} \mathbf{v}_{\mathcal{A}}}{\mathbf{v}'_{\mathcal{A}} \mathbf{v}_{\mathcal{A}}} \|\mathbf{\Delta}^b\|_2^2 \\ &\leq \max_{|\mathcal{A}| \lesssim \tilde{N}} \|\mathbf{E}_{\mathcal{A}}\|_2^2 \|\mathbf{\Delta}^b\|_2^2 \leq \max_{|\mathcal{A}| \lesssim \tilde{N}} \max_{\ell \in \{1, \dots, L_n\}} \|\tilde{\mathbf{E}}_{\mathcal{A}, \ell}\|_2^2 \left( \sum_{\ell=0}^{L_n} \|\Phi_{\ell}\|_2 \right)^2 \|\mathbf{\Delta}^b\|_2^2 \end{aligned}$$

where  $\mathbf{v}_{\mathcal{A}} \in \mathbb{R}^{|\mathcal{A}|}$  consists of elements  $\{v_i : i \in \mathcal{A}\}$  and  $\mathbf{E}_{\mathcal{A}} \in \mathbb{R}^{T \times |\mathcal{A}|}$  is composed of the corresponding columns. Note that the second inequality holds since  $\|\Delta^b \mathbf{u}\|_0 \lesssim \tilde{N}$ , and in the last inequality  $\tilde{\mathbf{E}}_{\mathcal{A},\ell}$  is defined in the proof of Lemma 1. We also observe that  $\sum_{\ell=0}^{\infty} \|\Phi_{\ell}\|_2 < \infty$  by Assumption 3. By Theorem 5.39 of Vershynin (2012) with the union bound, for some constants  $c_1$  and  $c_2$  such that  $c_1 < c_2$  and  $C$ , we have

$$\begin{aligned} & \mathbb{P} \left( \max_{|\mathcal{A}| \lesssim \tilde{N}} \max_{\ell \in \{0, \dots, L_n\}} \|\tilde{\mathbf{E}}_{\mathcal{A},\ell}\|_2 > C(\tilde{N} \vee T)^{1/2} \log^{1/2}(N \vee T) \right) \\ & \leq \binom{N}{c_1 \tilde{N}} (L_n + 1) \max_{|\mathcal{A}| \lesssim \tilde{N}} \max_{\ell \in \{1, \dots, L_n\}} \mathbb{P} \left( \|\tilde{\mathbf{E}}_{\mathcal{A},\ell}\|_2 > C(\tilde{N} \vee T)^{1/2} \log^{1/2}(N \vee T) \right) \\ & \lesssim N^{c_1 \tilde{N}} (N \vee T)^{\nu} \exp \left\{ -c_2 (\tilde{N} \vee T) \log(N \vee T) \right\} \\ & = O \left( (N \vee T)^{-\tilde{N} \vee T} \right) = O \left( (N \vee T)^{-\nu} \right). \end{aligned}$$

Thus, we have with probability at least  $1 - O((N \vee T)^{-\nu})$ ,

$$\begin{aligned} \|\mathbf{E} \Delta^b\|_2 & \lesssim \|\Delta^b\|_2 (\tilde{N} \vee T)^{1/2} \log^{1/2}(N \vee T) \\ & \leq \|\Delta^b\|_{\text{F}} (\tilde{N} \vee T)^{1/2} \log^{1/2}(N \vee T), \end{aligned}$$

giving the desired bound.  $\square$

**Lemma 5.** *Suppose the same conditions as Theorem 2. Then we have*

$$\|\Delta\|_{\text{F}}^2 \gtrsim \kappa_n \left( \|\hat{\mathbf{F}} - \mathbf{F}^0\|_{\text{F}}^2 + \|\hat{\mathbf{B}} - \mathbf{B}^0\|_{\text{F}}^2 \right),$$

where

$$\kappa_n = \frac{N_r(N_r \wedge T)}{N_1}.$$

*Proof.* Recall the notation based on the SVD of  $\mathbf{C}^0$  and  $\hat{\mathbf{C}}$ :  $\mathbf{U}^0 = \mathbf{F}^0$ ,  $\mathbf{V}^0 \mathbf{D}^0 = \mathbf{B}^0$ ,  $\hat{\mathbf{U}} = \hat{\mathbf{F}}$ , and  $\hat{\mathbf{V}} \hat{\mathbf{D}} = \hat{\mathbf{B}}$ . To establish the statement, we derive the following two inequalities:

$$\begin{aligned} (a) \quad \|\Delta\|_{\text{F}}^2 & \gtrsim \frac{N_r^2}{N_1} \|\hat{\mathbf{U}} - \mathbf{U}^0\|_{\text{F}}^2, \\ (b) \quad \|\Delta\|_{\text{F}}^2 & \gtrsim \frac{T N_r}{N_1} \|\hat{\mathbf{D}} \hat{\mathbf{V}}' - \mathbf{D}^0 \mathbf{V}^{0'}\|_{\text{F}}^2. \end{aligned}$$

Using them, we can immediately obtain the result.

First we prove (a). We define matrices:  $\hat{\mathbf{U}}_* = T^{-1/2} \hat{\mathbf{U}}$ ,  $\hat{\mathbf{D}}_* = \hat{\mathbf{D}} \hat{\mathbf{N}}^{1/2}$ ,  $\hat{\mathbf{V}}_* = \hat{\mathbf{V}} \hat{\mathbf{N}}^{-1/2}$ ,  $\mathbf{U}_*^0 = T^{-1/2} \mathbf{U}^0$ ,  $\mathbf{D}_*^0 = \mathbf{D}^0 \mathbf{N}^{1/2}$ , and  $\mathbf{V}_*^0 = \mathbf{V}^0 \mathbf{N}^{-1/2}$ , where  $\hat{\mathbf{N}}$  is any p.d. diagonal matrix. Then, we can see that

$$T^{-1/2} \Delta = \hat{\mathbf{U}}_* \hat{\mathbf{D}}_* \hat{\mathbf{V}}_*' - \mathbf{U}_*^0 \mathbf{D}_*^0 \mathbf{V}_*^{0'} =: \Delta_*.$$

For this expression, we can apply the proof of Lemma 3 in Uematsu et al. (2019). That is, under Assumptions 1 and 2, we have

$$\begin{aligned}
\|\hat{\mathbf{U}}_* - \mathbf{U}_*^0\|_F^2 &= \sum_{k=1}^r \|\hat{\mathbf{u}}_{*k} - \mathbf{u}_{*k}^0\|_2^2 \leq \|\Delta_*\|_F^2 (cd_{*1}^2/\delta) \sum_{k=1}^r d_{*k}^{-4} \\
&= cd_1^2 N_1 \|\Delta_*\|_F^2 \sum_{k=1}^r \frac{1}{\delta d_k^4 N_k^2} \\
&= \|\Delta_*\|_F^2 (cd_1^2 N_1/\delta) \sum_{k=1}^r (d_k N_k^{1/2})^{-4} \lesssim \|\Delta_*\|_F^2 \frac{N_1}{N_r^2}.
\end{aligned}$$

Rewriting this inequality with the original scaling gives result (a).

Next, we prove (b). We begin with rewriting  $\Delta_*$  as

$$\hat{\mathbf{U}}_*(\hat{\mathbf{D}}_* \hat{\mathbf{V}}_*' - \mathbf{D}_*^0 \mathbf{V}_*^{0'}) = \Delta_* - (\hat{\mathbf{U}}_* - \mathbf{U}_*^0) \mathbf{D}_*^0 \mathbf{V}_*^{0'}.$$

The triangle inequality and unitary property of the Frobenius norm entail that

$$\|\hat{\mathbf{D}}_* \hat{\mathbf{V}}_*' - \mathbf{D}_*^0 \mathbf{V}_*^{0'}\|_F \leq \|\Delta_*\|_F + \|(\hat{\mathbf{U}}_* - \mathbf{U}_*^0) \mathbf{D}_*^0\|_F.$$

We can bound the second term of the upper bound as in the proof of (a). That is, we have

$$\begin{aligned}
\|(\hat{\mathbf{U}}_* - \mathbf{U}_*^0) \mathbf{D}_*^0\|_F^2 &\leq \|\Delta_*\|_F^2 (cd_{*1}^2/\delta) \sum_{k=1}^r d_{*k}^{-2} \\
&= \|\Delta_*\|_F^2 (cd_1^2 N_1/\delta) \sum_{k=1}^r (d_k N_k^{1/2})^{-2} \lesssim \|\Delta_*\|_F^2 \frac{N_1}{N_r}.
\end{aligned}$$

Combining these inequalities gives

$$\begin{aligned}
\|\hat{\mathbf{D}}_* \hat{\mathbf{V}}_*' - \mathbf{D}_*^0 \mathbf{V}_*^{0'}\|_F^2 &\leq 2\|\Delta_*\|_F^2 + 2\|(\hat{\mathbf{U}}_* - \mathbf{U}_*^0) \mathbf{D}_*^0\|_F^2 \\
&\lesssim \|\Delta_*\|_F^2 + \|\Delta_*\|_F^2 \frac{N_1}{N_r} = T^{-1} \|\Delta\|_F^2 \left(1 + \frac{N_1}{N_r}\right).
\end{aligned}$$

Noting that the left-hand side is equal to  $\|\hat{\mathbf{D}} \hat{\mathbf{V}}' - \mathbf{D}^0 \mathbf{V}^{0'}\|_F^2$ , we obtain

$$\begin{aligned}
\|\Delta\|_F^2 &\gtrsim T \left(1 + \frac{N_1}{N_r}\right)^{-1} \|\hat{\mathbf{D}} \hat{\mathbf{V}}' - \mathbf{D}^0 \mathbf{V}^{0'}\|_F^2 \\
&= \frac{TN_r}{N_1 + N_r} \|\hat{\mathbf{D}} \hat{\mathbf{V}}' - \mathbf{D}^0 \mathbf{V}^{0'}\|_F^2 \gtrsim \frac{TN_r}{N_1} \|\hat{\mathbf{D}} \hat{\mathbf{V}}' - \mathbf{D}^0 \mathbf{V}^{0'}\|_F^2.
\end{aligned}$$

This completes the proof.  $\square$

**Lemma 6.** Suppose that Assumptions 1–4 with  $\tilde{N} = N$  and conditions (9) and (10) hold. Then we have

$$\|\hat{\mathbf{B}}_{\text{PC}} - \mathbf{B}^0\|_{\max} \lesssim \frac{N^{1/2} \gamma_n(N)}{(N \vee T)^{1/2}}$$

with probability at least  $1 - O((N \vee T)^{-\nu})$ .

*Proof.* By the definition of the PC estimator under PC1 restriction, we have

$$\begin{aligned}
\widehat{\mathbf{B}}_{\text{PC}} &= T^{-1} \mathbf{X}' \widehat{\mathbf{F}}_{\text{PC}} \\
&= T^{-1} (\mathbf{B}^0 \mathbf{F}^{0'} + \mathbf{E}') \widehat{\mathbf{F}}_{\text{PC}} \\
&= T^{-1} (\mathbf{B}^0 \mathbf{F}^{0'} + \mathbf{E}') \mathbf{F}^0 + T^{-1} (\mathbf{B}^0 \mathbf{F}^{0'} + \mathbf{E}') (\widehat{\mathbf{F}}_{\text{PC}} - \mathbf{F}^0) \\
&= \mathbf{B}^0 + T^{-1} \mathbf{E}' \mathbf{F}^0 + T^{-1} \mathbf{B}^0 \mathbf{F}^{0'} (\widehat{\mathbf{F}}_{\text{PC}} - \mathbf{F}^0) + T^{-1} \mathbf{E}' (\widehat{\mathbf{F}}_{\text{PC}} - \mathbf{F}^0),
\end{aligned}$$

which together with the triangle inequality implies that

$$\begin{aligned}
&\|\widehat{\mathbf{B}}_{\text{PC}} - \mathbf{B}^0\|_{\max} \\
&\leq T^{-1} \|\mathbf{E}' \mathbf{F}^0\|_{\max} + T^{-1} \|\mathbf{B}^0 \mathbf{F}^{0'} (\widehat{\mathbf{F}}_{\text{PC}} - \mathbf{F}^0)\|_{\max} + T^{-1} \|\mathbf{E}' (\widehat{\mathbf{F}}_{\text{PC}} - \mathbf{F}^0)\|_{\max}. \quad (\text{A.16})
\end{aligned}$$

From Lemma 1(c), the first term of (A.16) is bounded by  $T^{-1/2} \log^{1/2}(N \vee T)$  (up to a positive constant factor) with probability at least  $1 - O((N \vee T)^{-\nu})$ . We then consider the remaining two terms. The second term of (A.16) is evaluated as

$$\begin{aligned}
T^{-1} \|\mathbf{B}^0 \mathbf{F}^{0'} (\widehat{\mathbf{F}}_{\text{PC}} - \mathbf{F}^0)\|_{\max} &\leq r T^{-1} \|\mathbf{B}^0\|_{\max} \|\mathbf{F}^{0'} (\widehat{\mathbf{F}}_{\text{PC}} - \mathbf{F}^0)\|_{\max} \\
&\lesssim T^{-1} \|\mathbf{F}^0\|_{\text{F}} \|\widehat{\mathbf{F}}_{\text{PC}} - \mathbf{F}^0\|_{\text{F}} \\
&\lesssim T^{-1/2} \|\widehat{\mathbf{F}}_{\text{PC}} - \mathbf{F}^0\|_{\text{F}}.
\end{aligned}$$

By the Cauchy-Schwarz inequality, the third term of (A.16) is evaluated as

$$\begin{aligned}
T^{-1} \|\mathbf{E}' (\widehat{\mathbf{F}}_{\text{PC}} - \mathbf{F}^0)\|_{\max} &= T^{-1} \max_{i,k} \left| \sum_t e_{ti} (\hat{f}_{tk} - f_{tk}^0) \right| \\
&\leq \max_i \left( T^{-1} \sum_t e_{ti}^2 \right)^{1/2} \max_k \left( T^{-1} \sum_t (\hat{f}_{tk} - f_{tk}^0)^2 \right)^{1/2} \\
&\leq \max_i \left\{ \left| T^{-1} \sum_t (e_{ti}^2 - \mathbb{E} e_{ti}^2) \right|^{1/2} + \left| T^{-1} \sum_t \mathbb{E} e_{ti}^2 \right|^{1/2} \right\} T^{-1/2} \|\widehat{\mathbf{F}}_{\text{PC}} - \mathbf{F}^0\|_{\text{F}} \\
&\lesssim T^{-1/2} \|\widehat{\mathbf{F}}_{\text{PC}} - \mathbf{F}^0\|_{\text{F}},
\end{aligned}$$

where the last inequality follows from Lemma 1(d) and

$$\begin{aligned}
\|\mathbb{E} \mathbf{e}_t \mathbf{e}_t'\|_{\max} &\leq \sum_{\ell=0}^{L_n} \|\boldsymbol{\Phi}_\ell \mathbb{E} \boldsymbol{\varepsilon}_{t-\ell} \boldsymbol{\varepsilon}_{t-\ell}' \boldsymbol{\Phi}_\ell'\|_{\max} \\
&\lesssim \sum_{\ell=0}^{L_n} \|\boldsymbol{\Phi}_\ell \boldsymbol{\Phi}_\ell'\|_{\max} \lesssim \sum_{\ell=0}^{\infty} \|\boldsymbol{\Phi}_\ell\|_2^2 < \infty.
\end{aligned}$$

Consequently, by Theorem 3 with condition (10), the bound in (A.16) becomes

$$\begin{aligned}
\|\widehat{\mathbf{B}}_{\text{PC}} - \mathbf{B}^0\|_{\max} &\lesssim T^{-1/2} \log^{1/2}(N \vee T) + T^{-1/2} \|\widehat{\mathbf{F}}_{\text{PC}} - \mathbf{F}^0\|_{\text{F}} \\
&\lesssim \frac{\log^{1/2}(N \vee T)}{T^{1/2}} + \frac{N^{1/2} \gamma_n(N)}{(N \vee T)^{1/2}} \\
&= \frac{\log^{1/2}(N \vee T)}{T^{1/2}} + \frac{(N \wedge T)^{1/2} \gamma_n(N)}{T^{1/2}}
\end{aligned}$$



with probability at least  $1 - O((N \vee T)^{-\nu})$ . In this upper bound, the second term dominates the first term because

$$\begin{aligned} (N \wedge T)^{1/2} \gamma_n(N) &= \frac{(N \wedge T)^{1/2} N_1 (N \vee T)^{1/2} \log^{1/2}(N \vee T)}{N_r (N_r \wedge T)} \\ &= \frac{N_1}{N_r} \cdot \frac{N^{1/2} T^{1/2} \log^{1/2}(N \vee T)}{N_r \wedge T} \geq 1 \cdot \log^{1/2}(N \vee T). \end{aligned}$$

Thus, we have

$$\|\widehat{\mathbf{B}}_{\text{PC}} - \mathbf{B}^0\|_{\max} \asymp T^{-1/2} \|\widehat{\mathbf{F}}_{\text{PC}} - \mathbf{F}^0\|_{\text{F}} \lesssim \frac{N^{1/2} \gamma_n(N)}{(N \vee T)^{1/2}}$$

with probability at least  $1 - O((N \vee T)^{-\nu})$ . This completes the proof of Lemma 6.  $\square$

**Lemma 7.** *Suppose the same conditions as Theorem 4. Then, for any deterministic matrices  $\mathbf{U} = (u_{tk}) \in \mathbb{R}^{T \times r}$  and  $\mathbf{V} = (v_{ik}) \in \mathbb{R}^{N \times r}$ , the following inequalities simultaneously hold with probability at least  $1 - O((N \vee T)^{-\nu})$ :*

- (a)  $|\text{tr} \mathbf{E} \mathbf{B}^0 \mathbf{U}'| \lesssim N_1^{1/2} \|\mathbf{U}\|_{\text{F}} \log^{1/2}(N \vee T),$
- (b)  $|\text{tr} \mathbf{E}' \mathbf{F}^0 \mathbf{V}'_{\mathcal{S}}| \lesssim T^{1/2} \|\mathbf{V}_{\mathcal{S}}\|_{\text{F}} \log^{1/2}(N \vee T),$
- (c)  $|\text{tr} \mathbf{V}'_{\mathcal{S}} \mathbf{E}' \mathbf{U}| \lesssim \|\mathbf{U}\|_{\text{F}} \|\mathbf{V}_{\mathcal{S}}\|_{\text{F}} \log^{1/2}(N \vee T),$
- (d)  $|\text{tr} \mathbf{V}_{\mathcal{S}} \mathbf{U}' \mathbf{F}^0 \mathbf{V}'_{\mathcal{S}}| \lesssim \|\mathbf{U}\|_{\text{F}} \|\mathbf{V}_{\mathcal{S}}\|_{\text{F}}^2 \log^{1/2}(N \vee T),$
- (e)  $|\text{tr} \mathbf{B}^0 \mathbf{U}' \mathbf{U} \mathbf{V}'_{\mathcal{S}}| \lesssim N_1^{1/2} \|\mathbf{U}\|_{\text{F}}^2 \|\mathbf{V}_{\mathcal{S}}\|_{\text{F}},$
- (f)  $|\text{tr} \mathbf{B}^0 \mathbf{U}' \mathbf{F}^0 \mathbf{V}'_{\mathcal{S}}| \lesssim N_1^{1/2} \|\mathbf{U}\|_{\text{F}} \|\mathbf{V}_{\mathcal{S}}\|_{\text{F}} \log^{1/2}(N \vee T).$

*Proof.* Recall that  $\mathbf{V}_{\mathcal{S}} \in \mathbb{R}^{N \times r}$  is defined as the matrix whose  $(i, k)$ th element is  $v_{ik} 1\{(i, k) \in \mathcal{S}\}$ , where  $\mathcal{S} = \text{supp}(\mathbf{B}^0)$ ; see the proof of Theorem 4.

(a) First note that the  $(t, k)$ th element of  $\mathbf{E} \mathbf{B}^0$  is given by  $\mathbf{e}'_t \mathbf{b}_k^0$ . We observe that

$$|\text{tr} \mathbf{E} \mathbf{B}^0 \mathbf{U}'| = |\text{vec}(\mathbf{E} \mathbf{B}^0)' \mathbf{u}| \leq r \max_k \left| \sum_{t=1}^T \mathbf{e}'_t \mathbf{b}_k^0 u_{tk} \right|,$$

where we have written as  $\mathbf{u} = \text{vec}(\mathbf{U})$ . From Assumption 3, recall that  $\mathbf{e}_t = \sum_{\ell=0}^L \Phi_{\ell} \boldsymbol{\varepsilon}_{t-\ell}$ , where  $\boldsymbol{\varepsilon}_t = (\varepsilon_{t1}, \dots, \varepsilon_{tN})'$  with  $\{\varepsilon_{ti}\}_{t,i} \sim \text{i.i.d. subG}(\sigma_{\varepsilon}^2)$ . Let  $\tilde{b}_{\ell k, i}$  denote the  $i$ th element of

$\Phi'_\ell \mathbf{b}_k^0$  as in the proof of Lemma 1(b). Then, we have

$$\begin{aligned}
\max_k \left| \sum_{t=1}^T \mathbf{e}'_t \mathbf{b}_k^0 u_{tk} \right| &= \max_k \left| \sum_{t=1}^T \sum_{\ell=0}^L \sum_{i=1}^N \varepsilon_{t-\ell, i} \tilde{b}_{\ell k, i} u_{tk} \right| \\
&\leq \sum_{\ell=0}^L \max_k \left| \sum_{t=1}^T \sum_{i=1}^N \varepsilon_{t-\ell, i} \tilde{b}_{\ell k, i} u_{tk} \right| \\
&\leq \sum_{\ell=0}^L \max_k \left\| \Phi'_\ell \mathbf{b}_k \right\|_2^{-1} \sum_{t=1}^T \sum_{i=1}^N \varepsilon_{t-\ell, i} \tilde{b}_{\ell k, i} u_{tk} \left\| \Phi'_\ell \mathbf{b}_k \right\|_2 \\
&\leq \max_{k, \ell} \left\| \Phi'_\ell \mathbf{b}_k \right\|_2^{-1} \sum_{t=1}^T \sum_{i=1}^N \varepsilon_{t-\ell, i} \tilde{b}_{\ell k, i} u_{tk} \max_k \left\| \mathbf{b}_k \right\|_2 \sum_{\ell=0}^\infty \left\| \Phi_\ell \right\|_2 \\
&\lesssim N_1^{1/2} \max_{k, \ell} \left| \sum_{t=1}^T u_{tk} \left\| \Phi'_\ell \mathbf{b}_k \right\|_2^{-1} \sum_{i=1}^N \varepsilon_{t-\ell, i} \tilde{b}_{\ell k, i} \right|.
\end{aligned}$$

Since  $\{\varepsilon_{t-\ell, i} \tilde{b}_{\ell k, i}\}_i$  is a sequence of i.i.d.  $\text{subG}(\sigma_\varepsilon^2 \tilde{b}_{\ell k, i}^2)$  for each  $t, k, \ell$ , we can see that  $\{\left\| \Phi'_\ell \mathbf{b}_k \right\|_2^{-1} \sum_{i=1}^N \varepsilon_{t-\ell, i} \tilde{b}_{\ell k, i}\}_t \sim \text{i.i.d. subG}(\sigma_\varepsilon^2)$  by Lemma 2. Moreover, Lemma 2 gives

$$Z_{k\ell} := \sum_{t=1}^T u_{tk} \left\| \Phi'_\ell \mathbf{b}_k \right\|_2^{-1} \sum_{i=1}^N \varepsilon_{t-\ell, i} \tilde{b}_{\ell k, i} \sim \text{subG}(\sigma_\varepsilon^2 \|\mathbf{u}_k\|_2^2).$$

Therefore, the subG tail inequality and the union bound entail

$$\begin{aligned}
\mathbb{P} \left( \max_{k, \ell} |Z_{k\ell}| > x \right) &\leq r(L+1) \max_{k, \ell} \mathbb{P}(|Z_{k\ell}| > x) \\
&\leq 2r(N \vee T)^\nu \exp \left( -\frac{x^2}{2\sigma_\varepsilon^2 \max_k \|\mathbf{u}_k\|_2^2} \right) \\
&\leq 2r(N \vee T)^\nu \exp \left( -\frac{x^2}{2\sigma_\varepsilon^2 \|\mathbf{U}\|_F^2} \right).
\end{aligned}$$

Setting  $x^2 = 4\sigma_\varepsilon^2 \|\mathbf{U}\|_F^2 \nu \log(N \vee T)$  leads to getting the bound

$$\max_{k, \ell} |Z_{k\ell}| \leq 2\sigma_\varepsilon \|\mathbf{U}\|_F \nu^{1/2} \log^{1/2}(N \vee T)$$

with probability at least  $1 - O((N \vee T)^{-\nu})$ . Thus the desired upper bound

$$|\text{tr} \mathbf{E} \mathbf{B}^0 \mathbf{U}'| \lesssim N_1^{1/2} \log^{1/2}(N \vee T) \|\mathbf{U}\|_F$$

holds with probability at least  $1 - O((N \vee T)^{-\nu})$ .

(b) As in the proof of Lemma 1, we write  $\tilde{\mathbf{E}}_\ell = (\varepsilon_{1-\ell}, \dots, \varepsilon_{T-\ell})' \in \mathbb{R}^{T \times N}$  and  $\tilde{\mathbf{Z}}_\ell = (\zeta_{1-\ell}, \dots, \zeta_{T-\ell})' \in \mathbb{R}^{T \times r}$ . Then we can write  $\mathbf{E}' \mathbf{F} = \sum_{\ell, m=0}^{L_n} \Phi'_\ell \tilde{\mathbf{E}}'_\ell \tilde{\mathbf{Z}}_m \Psi'_m$  under Assumptions 1 and 3. By the same way as in (a), we have

$$\begin{aligned}
|\text{tr} \mathbf{E}' \mathbf{F}^0 \mathbf{V}'_S| &= \left| \sum_{(i, k) \in S} \sum_{\ell, m=0}^{L_n} \phi'_{\ell, i} \tilde{\mathbf{E}}'_\ell \tilde{\mathbf{Z}}_m \psi_{m, k} v_{ik} \right| \leq \sum_{\ell, m=0}^{L_n} \left| \sum_{(i, k) \in S} \phi'_{\ell, i} \tilde{\mathbf{E}}'_\ell \tilde{\mathbf{Z}}_m \psi_{m, k} v_{ik} \right| \\
&= \sum_{\ell, m=0}^{L_n} \left| \sum_{(i, k) \in S} v_{ik} \text{tr} \psi_{m, k} \phi'_{\ell, i} \tilde{\mathbf{E}}'_\ell \tilde{\mathbf{Z}}_m \right| = \sum_{\ell, m=0}^{L_n} \left| \text{tr} \Theta_{\ell m} \tilde{\mathbf{E}}'_\ell \tilde{\mathbf{Z}}_m \right|,
\end{aligned}$$

where  $\Theta_{\ell m} := \sum_{(i,k) \in \mathcal{S}} v_{ik} \psi_{m,k} \phi'_{\ell,i}$  with its  $(h, j)$ th component given by  $\theta_{\ell m, hj}$  for  $h = 1, \dots, r$  and  $j = 1, \dots, N$ . Recall that  $\tilde{\mathbf{E}}'_\ell = (\varepsilon_{1-\ell}, \dots, \varepsilon_{T-\ell})$  and  $\tilde{\mathbf{Z}}'_m = (\zeta_{1-m}, \dots, \zeta_{T-m})$  from the proof of Lemma 1. Then we have

$$\begin{aligned} \sum_{\ell, m=0}^{L_n} \left| \text{tr } \Theta_{\ell m} \tilde{\mathbf{E}}'_\ell \tilde{\mathbf{Z}}'_m \right| &= \sum_{\ell, m=0}^{L_n} \left| \sum_{h=1}^r \sum_{t=1}^T \left( \sum_{j=1}^N \theta_{\ell m, hj} \varepsilon_{t-\ell, j} \right) \zeta_{t-m, h} \right| \\ &\leq r \max_h \sum_{\ell, m=0}^{L_n} \left| \sum_{t=1}^T \left( \|\boldsymbol{\theta}_{\ell m, h}\|_2^{-1} \sum_{j=1}^N \theta_{\ell m, hj} \varepsilon_{t-\ell, j} \right) \zeta_{t-m, h} \right| \|\boldsymbol{\theta}_{\ell m, h}\|_2 \\ &\lesssim \max_{h, \ell, m} \left| \sum_{t=1}^T \left( \|\boldsymbol{\theta}_{\ell m, h}\|_2^{-1} \sum_{j=1}^N \theta_{\ell m, hj} \varepsilon_{t-\ell, j} \right) \zeta_{t-m, h} \right| \sum_{\ell, m=0}^{L_n} \|\boldsymbol{\theta}_{\ell m, h}\|_2 \\ &\lesssim \max_{h, \ell, m} \left| \sum_{t=1}^T \left( \|\boldsymbol{\theta}_{\ell m, h}\|_2^{-1} \sum_{j=1}^N \theta_{\ell m, hj} \varepsilon_{t-\ell, j} \right) \zeta_{t-m, h} \right| \max_h \sum_{\ell, m=0}^{L_n} \|\boldsymbol{\theta}_{\ell m, h}\|_2, \end{aligned}$$

where  $\boldsymbol{\theta}'_{\ell m, h}$  is the  $h$ th row vector of  $\Theta_{\ell m}$ . By the same reason as in the proof of Lemma 1(c), Lemma 2 entails that the inside of the absolute value is the sum of i.i.d.  $\text{subE}(4e\sigma_\varepsilon\sigma_\zeta)$  random variables. Thus, the same bound in that proof can be used. Thus, applying the union bound, we obtain with probability at least  $1 - O((N \vee T)^{-\nu})$ ,

$$\max_{h, \ell, m} \left| \sum_{t=1}^T \left( \|\boldsymbol{\theta}_{\ell m, h}\|_2^{-1} \sum_{j=1}^N \theta_{\ell m, hj} \varepsilon_{t-\ell, j} \right) \zeta_{t-m, h} \right| \leq (96e^2\sigma_\varepsilon^2\sigma_\zeta^2\nu T \log(N \vee T))^{1/2}.$$

Finally, we evaluate  $\max_h \sum_{\ell, m=0}^{L_n} \|\boldsymbol{\theta}_{\ell m, h}\|_2$ . By the construction, we have

$$\begin{aligned} \max_h \sum_{\ell, m=0}^{L_n} \|\boldsymbol{\theta}_{\ell m, h}\|_2 &= \max_h \sum_{\ell, m=0}^{L_n} \left( \sum_{j=1}^N \left( \sum_{(i,k) \in \mathcal{S}} v_{ik} \psi_{m, hk} \phi_{\ell, ij} \right)^2 \right)^{1/2} \\ &\leq \max_h \sum_{\ell, m=0}^{L_n} \left( \sum_{k=1}^r \psi_{m, hk}^2 \sum_{i,j=1}^N \phi_{\ell, ij}^2 \right)^{1/2} \|\mathbf{v}_S\|_2 \\ &\leq \sum_{m=0}^{\infty} \|\Psi_m\|_2 \sum_{\ell=0}^{\infty} \|\Phi_\ell\|_F \|\mathbf{V}_S\|_F \lesssim \|\mathbf{V}_S\|_F. \end{aligned}$$

Thus the desired upper bound

$$|\text{tr } \mathbf{E}' \mathbf{F}^0 \mathbf{V}'_S| \lesssim T^{1/2} \log^{1/2}(N \vee T) \|\mathbf{V}_S\|_F$$

holds with probability at least  $1 - O((N \vee T)^{-\nu})$ .

(c) We observe that

$$|\text{tr } \mathbf{V}'_S \mathbf{E}' \mathbf{U}| = \left| \sum_{k=1}^r \sum_{t=1}^T \mathbf{v}'_k \mathbf{e}_t u_{tk} \right| \leq \sum_{k=1}^r \sum_{\ell=0}^L \left| \sum_{t=1}^T \mathbf{v}'_k \Phi_\ell \varepsilon_{t-\ell} u_{tk} \right|.$$

By Assumption 3 and Lemma 2, we have  $(\mathbf{v}'_k \boldsymbol{\Phi}_\ell \boldsymbol{\varepsilon}_{t-\ell})_t \sim \text{indep. subG}(\sigma_\varepsilon^2 \|\mathbf{v}'_k \boldsymbol{\Phi}_\ell\|_2^2)$  for each  $k$  and  $\ell$ . Thus, by Lemma 2 again, we further have  $\sum_{t=1}^T \mathbf{v}'_k \boldsymbol{\Phi}_\ell \boldsymbol{\varepsilon}_{t-\ell} u_{tk} \sim \text{subG}(\sigma_\varepsilon^2 \|\mathbf{v}'_k \boldsymbol{\Phi}_\ell\|_2^2 \|\mathbf{u}_k\|_2^2)$  for each  $k$  and  $\ell$ . Therefore, the subG tail probability gives

$$\begin{aligned} \left| \sum_{t=1}^T \mathbf{v}'_k \boldsymbol{\Phi}_\ell \boldsymbol{\varepsilon}_{t-\ell} u_{tk} \right| &\lesssim \|\mathbf{v}'_k \boldsymbol{\Phi}_\ell\|_2 \|\mathbf{u}_k\|_2 \log^{1/2}(N \vee T) \\ &\leq \|\boldsymbol{\Phi}_\ell\|_2 \|\mathbf{V}_S\|_F \|\mathbf{U}\|_F \log^{1/2}(N \vee T) \end{aligned}$$

with probability at least  $1 - O((N \vee T)^{-\nu})$ . Consequently, we have

$$\begin{aligned} |\text{tr } \mathbf{V}'_S \mathbf{E}' \mathbf{U}| &\lesssim \sum_{\ell=0}^{\infty} \|\boldsymbol{\Phi}_\ell\|_2 \|\mathbf{V}_S\|_F \|\mathbf{U}\|_F \log^{1/2}(N \vee T) \\ &\lesssim \|\mathbf{V}_S\|_F \|\mathbf{U}\|_F \log^{1/2}(N \vee T), \end{aligned}$$

which yields the result.

(d) By the property of norms, we obtain

$$\begin{aligned} |\text{tr } \mathbf{V}'_S \mathbf{V}_S \mathbf{U}' \mathbf{F}^0| &\leq \|\mathbf{V}'_S \mathbf{V}_S\|_* \|\mathbf{U}' \mathbf{F}^0\|_2 \\ &\leq r^{3/2} \|\mathbf{V}'_S \mathbf{V}_S\|_F \|\mathbf{U}' \mathbf{F}^0\|_{\max} \lesssim \|\mathbf{V}_S\|_F^2 \max_{j,k} \left| \sum_{t=1}^T u_{tj} f_{tk}^0 \right|. \end{aligned}$$

By Assumption 1, the last stochastic part is evaluated as

$$\begin{aligned} \max_{j,k} \left| \sum_{t=1}^T u_{tk} f_{tk}^0 \right| &= \max_{j,k} \left| \sum_{\ell=0}^{L_n} \sum_{m=1}^r \psi_{\ell,km} \sum_{t=1}^T u_{tj} \zeta_{t-\ell,m} \right| \\ &\leq r \max_{k,m} \sum_{\ell=0}^{L_n} |\psi_{\ell,km}| \max_{j,m} \left| \sum_{t=1}^T \zeta_{t-\ell,m} u_{tj} \right| \\ &\leq r \max_{j,m,\ell} \left| \sum_{t=1}^T \zeta_{t-\ell,m} u_{tj} \right| \max_{k,m} \sum_{\ell=0}^{L_n} |\psi_{\ell,km}| \\ &\lesssim \max_{j,m,\ell} \left| \sum_{t=1}^T \zeta_{t-\ell,m} u_{tj} \right| \sum_{\ell=0}^{\infty} \|\boldsymbol{\Psi}_\ell\|_2, \end{aligned}$$

where  $\{\zeta_{tm}\}_{t,m} \sim \text{i.i.d. subG}(\sigma_\zeta^2)$  and  $\sum_{\ell=0}^{\infty} \|\boldsymbol{\Psi}_\ell\|_2$  is bounded. By Lemma 2(b), we have  $\sum_{t=1}^T \zeta_{t-\ell,m} u_{tj} \sim \text{subG}(\sigma_\zeta^2 \|\mathbf{u}_j\|_2^2)$  for any  $j, m, \ell$ . Thus, the subG tail inequality together with the union bound establishes that

$$\begin{aligned} \mathbb{P} \left( \max_{j,m,\ell} \left| \sum_{t=1}^T \zeta_{t-\ell,m} u_{tj} \right| > x \right) &\leq r^2 (L_n + 1) \max_{j,m,\ell} \mathbb{P} \left( \left| \sum_{t=1}^T \zeta_{t-\ell,m} u_{tj} \right| > x \right) \\ &\lesssim (N \vee T)^\nu \exp \left( -\frac{x^2}{2\sigma_\zeta^2 \max_j \|\mathbf{u}_j\|_2^2} \right). \end{aligned}$$

Setting  $x = 2\nu^{1/2} \sigma_\zeta \max_j \|\mathbf{u}_j\|_2 \log^{1/2}(N \vee T)$  yields

$$\begin{aligned} \max_{j,m,\ell} \left| \sum_{t=1}^T \zeta_{t-\ell,m} u_{tj} \right| &\leq 2\sigma_\zeta \max_j \|\mathbf{u}_j\|_2 \log^{1/2}(N \vee T) \\ &\lesssim \|\mathbf{U}\|_F \log^{1/2}(N \vee T) \end{aligned}$$

with probability at least  $1 - O((N \vee T)^{-\nu})$ . This together with the first inequality yields the result.

(e) We observe that

$$|\text{tr } \mathbf{B}^0 \mathbf{U}' \mathbf{U} \mathbf{V}'_{\mathcal{S}}| \leq \|\mathbf{V}'_{\mathcal{S}} \mathbf{B}^0\|_{\text{F}} \|\mathbf{U}' \mathbf{U}\|_{\text{F}} \lesssim N_1^{1/2} \|\mathbf{U}\|_{\text{F}}^2 \|\mathbf{V}_{\mathcal{S}}\|_{\text{F}},$$

which gives the proof.

(f) By the property of norms, we obtain

$$\begin{aligned} |\text{tr } \mathbf{V}'_{\mathcal{S}} \mathbf{B}^0 \mathbf{U}' \mathbf{F}^0| &\leq \|\mathbf{V}'_{\mathcal{S}} \mathbf{B}^0\|_* \|\mathbf{U}' \mathbf{F}^0\|_2 \\ &\leq r^{3/2} \|\mathbf{V}'_{\mathcal{S}} \mathbf{B}^0\|_{\text{F}} \|\mathbf{U}' \mathbf{F}^0\|_{\max} \lesssim N_1^{1/2} \|\mathbf{V}_{\mathcal{S}}\|_{\text{F}} \max_{j,k} \left| \sum_{t=1}^T u_{tj} f_{tk}^0 \right|. \end{aligned}$$

Thus by the same argument as the proof of (d), we conclude that the stochastic part is bounded by  $\|\mathbf{U}\|_{\text{F}} \log^{1/2}(N \vee T)$ , which occurs with probability at least  $1 - O((N \vee T)^{-\nu})$ . This completes the proofs of (a)–(f).  $\square$

**Lemma 8.** *Suppose the same conditions as Theorem 4. Then we have*

$$\|\mathbf{W}_{\mathcal{S}}\|_{\text{F}} \leq \frac{2(rN_1)^{1/2}}{\underline{b}_n^0}$$

with probability at least  $1 - O((N \vee T)^{-\nu})$ .

*Proof.* Let  $\hat{\underline{b}}_n = \min_{(i,k) \in \mathcal{S}} |\hat{b}_{ik}^{\text{ini}}|$ . For any  $x > 0$ , we have

$$\mathbb{P}(\|\mathbf{W}_{\mathcal{S}}\|_{\text{F}} > x) \leq \mathbb{P}\left(\|\mathbf{W}_{\mathcal{S}}\|_{\text{F}} > x \mid \hat{\underline{b}}_n > \underline{b}_n^0/2\right) + \mathbb{P}\left(\hat{\underline{b}}_n \leq \underline{b}_n^0/2\right). \quad (\text{A.17})$$

With setting  $x = 2(rN_1)^{1/2}/\underline{b}_n^0$ , we verify that the upper bound of (A.17) tends to zero. The first probability of the upper bound is bounded as

$$\begin{aligned} \mathbb{P}\left(\|\mathbf{W}_{\mathcal{S}}\|_{\text{F}} > \frac{2(rN_1)^{1/2}}{\underline{b}_n^0} \mid \hat{\underline{b}}_n > \underline{b}_n^0/2\right) &\leq \mathbb{P}\left(\frac{rN_1}{\hat{\underline{b}}_n^2} > \frac{4rN_1}{(\underline{b}_n^0)^2} \mid \hat{\underline{b}}_n > \underline{b}_n^0/2\right) \\ &\leq \mathbb{P}\left(\frac{2}{\hat{\underline{b}}_n \underline{b}_n^0} > \frac{4}{(\underline{b}_n^0)^2} \mid \hat{\underline{b}}_n > \underline{b}_n^0/2\right) \\ &= \mathbb{P}\left(\underline{b}_n^0/2 > \hat{\underline{b}}_n \mid \hat{\underline{b}}_n > \underline{b}_n^0/2\right) = 0. \end{aligned}$$

By condition (17) and Lemma 6, the second probability of the upper bound of (A.17) is bounded as

$$\begin{aligned} \mathbb{P}\left(\hat{\underline{b}}_n \leq \underline{b}_n^0/2\right) &\leq \mathbb{P}\left(\|\hat{\mathbf{B}}_{\text{ini}} - \mathbf{B}^0\|_{\max} \geq \underline{b}_n^0/2\right) \\ &\leq \mathbb{P}\left(\|\hat{\mathbf{B}}_{\text{ini}} - \mathbf{B}^0\|_{\max} \gtrsim \frac{N^{1/2} \gamma_n(\tilde{N})}{(\tilde{N} \vee T)^{1/2}}\right) = O((N \vee T)^{-\nu}). \end{aligned}$$

Consequently, we obtain

$$\|\mathbf{W}_{\mathcal{S}}\|_{\text{F}} \leq \frac{2(rN_1)^{1/2}}{\underline{b}_n^0}$$

with probability at least  $1 - O((N \vee T)^{-\nu})$ . This completes the proof.  $\square$

**Lemma 9.** *Suppose the same conditions as Theorem 4. Then we have*

$$\left\| \mathbf{W}_{\mathcal{S}^c}^- \circ (\mathbf{X}' \widehat{\mathbf{F}})_{\mathcal{S}^c} \right\|_{\max} < \eta_n$$

with probability at least  $1 - O((N \vee T)^{-\nu})$ .

*Proof.* Let  $\Delta = (\delta_{tk}) = \mathbf{F} - \mathbf{F}^0$  and  $\widehat{\Delta} = \widehat{\mathbf{F}} - \mathbf{F}^0$ . Define

$$\mathcal{F} = \{ \Delta \in \mathbb{R}^{T \times r} : \|\Delta\|_{\mathbf{F}} \leq Cr_n \}$$

where  $C$  is some positive constant introduced in the proof of Theorem 4 and

$$r_n = \frac{N_1(N_1 T)^{1/2} \log^{1/2}(N \vee T)}{N_r(N_r \wedge T)} = \frac{\gamma_n(\tilde{N}) N_1^{1/2} T^{1/2}}{(\tilde{N} \vee T)^{1/2}}.$$

Then we have

$$\begin{aligned} & \left\| \mathbf{W}_{\mathcal{S}^c}^- \circ (\mathbf{X}' \widehat{\mathbf{F}})_{\mathcal{S}^c} \right\|_{\max} \leq \left\| \mathbf{W}_{\mathcal{S}^c}^- \right\|_{\max} \left\| (\mathbf{X}' \widehat{\mathbf{F}})_{\mathcal{S}^c} \right\|_{\max} \\ & = \left\| \widehat{\mathbf{B}}_{\mathcal{S}^c}^{\text{ini}} \right\|_{\max} \left\| (\mathbf{B}^0 \mathbf{F}^{0'} \widehat{\Delta})_{\mathcal{S}^c} + (\mathbf{E}' \widehat{\Delta})_{\mathcal{S}^c} + (\mathbf{E}' \mathbf{F}^0)_{\mathcal{S}^c} \right\|_{\max} \\ & \leq \left\| \widehat{\mathbf{B}}^{\text{ini}} - \mathbf{B}^0 \right\|_{\max} \left( \sup_{\Delta \in \mathcal{F}} \left\| (\mathbf{B}^0 \mathbf{F}^{0'} \Delta)_{\mathcal{S}^c} \right\|_{\max} + \sup_{\Delta \in \mathcal{F}} \left\| (\mathbf{E}' \Delta)_{\mathcal{S}^c} \right\|_{\max} + \left\| (\mathbf{E}' \mathbf{F}^0)_{\mathcal{S}^c} \right\|_{\max} \right) \\ & \leq \left\| \widehat{\mathbf{B}}^{\text{ini}} - \mathbf{B}^0 \right\|_{\max} \left( \sup_{\Delta \in \mathcal{F}} \left\| \mathbf{B}^0 \mathbf{F}^{0'} \Delta \right\|_{\max} + \sup_{\Delta \in \mathcal{F}} \left\| \mathbf{E}' \Delta \right\|_{\max} + \left\| \mathbf{E}' \mathbf{F}^0 \right\|_{\max} \right). \end{aligned}$$

We evaluate the three terms in the parenthesis. Lemma 1(c) gives  $\|\mathbf{E}' \mathbf{F}^0\|_{\max} \lesssim T^{1/2} \log^{1/2}(N \vee T)$ . We next observe that for any  $\Delta$ ,

$$\left\| \mathbf{B}^0 \mathbf{F}^{0'} \Delta \right\|_{\max} \leq r \left\| \mathbf{B}^0 \right\|_{\max} \left\| \mathbf{F}^{0'} \Delta \right\|_{\max} \lesssim \left\| \mathbf{F}^{0'} \Delta \right\|_{\max}.$$

The upper bound is further bounded as

$$\left\| \mathbf{F}^{0'} \Delta \right\|_{\max} = \max_k \left\| \sum_t \mathbf{f}_t^0 \delta_{tk} \right\|_{\max} \leq \max_k \sum_{\ell} \left\| \Psi_{\ell} \sum_t \zeta_{t-\ell} \delta_{tk} \right\|_{\max}.$$

By Lemma 2 with Assumption 1, we have  $z_{\ell,jk} := \sum_t \zeta_{t-\ell,j} \delta_{tk} \sim \text{subG}(\sigma_{\zeta}^2 \|\delta_k\|_2^2)$  for each fixed  $\delta_{tk}$ ,  $j$ , and  $\ell$ . By the independence of  $z_{\ell,jk}$  across  $j$  and Lemma 2 again, we have  $\sum_j \psi_{\ell,ij} z_{\ell,jk} \sim \text{subG}(\sigma_{\zeta}^2 \|\delta_k\|_2^2 \|\Psi_{\ell,i}\|_2^2)$  for each  $i$ ,  $k$ , and  $\ell$ . Therefore, for any fixed  $\Delta$  and  $\ell$ , the subG tail inequality with the union bound entails that

$$\max_k \left\| \Psi_{\ell} \sum_t \zeta_{t-\ell} \delta_{tk} \right\|_{\max} \lesssim \max_i \|\Psi_{\ell,i}\|_2 \|\Delta\|_{\mathbf{F}} \log^{1/2}(N \vee T)$$

with probability at least  $1 - O((N \vee T)^{-\nu})$ . Because  $\max_i \|\Psi_{\ell,i}\|_2 \leq \|\Psi_{\ell}\|_2$  by the definition of the spectral norm, we have

$$\sup_{\Delta \in \mathcal{F}} \left\| \mathbf{F}^{0'} \Delta \right\|_{\max} \leq C \sum_{\ell} \|\Psi_{\ell}\|_2 r_n \log^{1/2}(N \vee T) \lesssim r_n \log^{1/2}(N \vee T)$$

with probability at least  $1 - O((N \vee T)^{-\nu})$ . Moreover, by the same argument as above with Assumption 3, we have

$$\sup_{\Delta \in \mathcal{F}} \|\mathbf{E}'\Delta\|_{\max} \lesssim r_n \log^{1/2}(N \vee T)$$

with probability at least  $1 - O((N \vee T)^{-\nu})$ . Combining the obtained bounds yields

$$\eta_n^{-1} \left\| \mathbf{W}_{\mathcal{S}^c}^- \circ (\mathbf{X}'\widehat{\mathbf{F}})_{\mathcal{S}^c} \right\|_{\max} \lesssim \eta_n^{-1} \left\| \widehat{\mathbf{B}}^{\text{ini}} - \mathbf{B}^0 \right\|_{\max} \left( T^{1/2} + r_n \right) \log^{1/2}(N \vee T).$$

We can obtain  $T^{1/2} + r_n = O(T^{1/2})$  because condition (10) entails

$$\frac{r_n}{T^{1/2}} = \frac{N_1^{3/2}}{N_r(N_r \wedge T)} = o\left(\frac{N_1^{3/2}}{N_1(N \vee T)^{1/2}}\right) = o\left(\frac{N_1^{1/2}}{(N \vee T)^{1/2}}\right) = o(1).$$

Hence, by Lemma 6 and condition (18), we have

$$\begin{aligned} \eta_n^{-1} \left\| \mathbf{W}_{\mathcal{S}^c}^- \circ (\mathbf{X}'\widehat{\mathbf{F}})_{\mathcal{S}^c} \right\|_{\max} &\lesssim \eta_n^{-1} \left\| \widehat{\mathbf{B}}^{\text{ini}} - \mathbf{B}^0 \right\|_{\max} T^{1/2} \log^{1/2}(N \vee T) \\ &\lesssim \frac{(\tilde{N} \vee T)^{1/2}}{N^{1/2} T^{1/2}} \frac{N^{1/2} \gamma_n(\tilde{N})}{(\tilde{N} \vee T)^{1/2}} T^{1/2} \log^{1/2}(N \vee T) \\ &= \gamma_n(\tilde{N}) \log^{1/2}(N \vee T) \end{aligned}$$

with probability at least  $1 - O((N \vee T)^{-\nu})$ . The last value is  $o(1)$  since  $\gamma_n(\tilde{N})$  tends to zero polynomially. This asymptotically guarantees the desired strict inequality.  $\square$

## C Additional Simulation Result: Global and Local Factors

In empirical data analysis, often the factor models are applied to capture the dependence within a sector or group, even hierarchical structure. For example, Ando and Bai (2017) consider two types of factors: global factors of which factor loadings have non-zero values for all the cross-section units, whilst the local factors of which loadings of specific cross sectional group are non-zero, and zero otherwise. Estimation of this class of hierarchical factor structure is of great interest in practice, and some estimation methods have been proposed.

Our WF-SOFAR estimator is easily applicable to estimate the model. To see how useful our estimator, we implemented a small experiment. We generate the data of four factors models,  $x_{ti} = \sum_{k=1}^r b_{ik} f_{tk} + e_{ti}$ , where  $f_{tk} = \rho_{fk} f_{t-1,k} + v_{tk}$  with  $v_{tk} \sim \text{i.i.d.} N(0, 1 - \rho_{fk}^2)$  and  $f_{1k} \sim \text{i.i.d.} N(0, 1)$ ;  $e_{ti} = \rho_e e_{t-1,i} + \beta \varepsilon_{t,i-1} + \beta \varepsilon_{t,i+1} + \varepsilon_{ti} \varepsilon_{ti} \sim \text{i.i.d.} N(0, \sigma_{\varepsilon,ti}^2)$  and  $\sigma_{\varepsilon,ti}^2$  is set so that  $\text{Var}(e_{ti}) = 1$  for  $i = 1, \dots, N$ ,  $t = 1, \dots, T$ . We set  $r = 4$ . The first factor is a global factor, i.e.,  $b_{i1} \sim \text{i.i.d.} N(0, 1)$  for  $i = 1, \dots, N$ . The other three factors are local ones, i.e.,  $b_{i2}$  is drawn from  $N(0, 1)$  for the first third,  $b_{i3}$  for the second third, and  $b_{i4}$  for the last third of cross section units while the rests are zero. We obtained a simulated data with  $N = 450$  and  $T = 120$ , and estimated the factor model given  $r = 4$  by the PC and WF-SOFAR. To visualize the quality of the factor loadings, we provide heat maps of three  $N \times N$  matrices,  $\sum_{k=1}^4 \omega_k \text{abs}(\mathbf{b}_k^0 \mathbf{b}_k^{0'})$ ,  $\sum_{k=1}^4 \omega_k \text{abs}(\widehat{\mathbf{b}}_k \widehat{\mathbf{b}}_k')$  and  $\sum_{k=1}^4 \omega_k \text{abs}(\widehat{\mathbf{b}}_{\text{PC},k} \widehat{\mathbf{b}}_{\text{PC},k}')$ . To clarify the difference between the global factor loadings and local ones, which overlaps, we use the weight  $\omega_1 = 1/8$  and  $\omega_2 = \omega_3 = \omega_4 = 1$ .

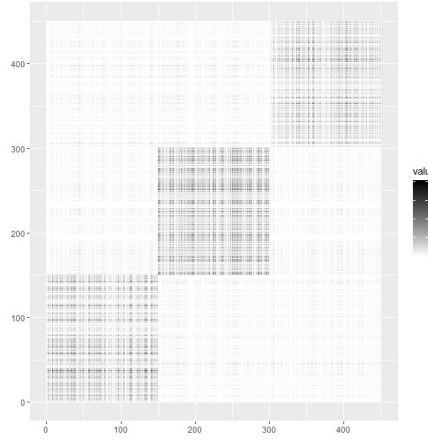


Figure 2: True factor loadings

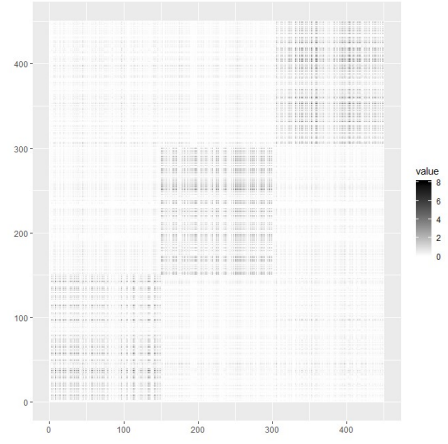


Figure 3: WF-SOFAR estimate

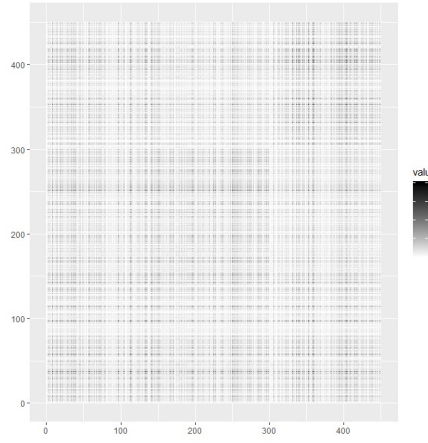


Figure 4: PC estimate



## D Additional Results of Empirical Example 1: Firm Security Returns

In addition to reporting the divergence rates, we summarize the estimates of the factor loadings, focusing on analysis of the contributions of industrial sectors to the non-zero factor loadings. Such contributions can be regarded as measures of sensitivities of industrial sectors to the factor. Also we look into the signs of the factor loadings. Notice that the firm securities with negative loadings react to the factor in the opposite direction to those with positive loadings. Therefore, given the systematic risk factor, the different sign of the factor loadings could be interpreted as the different investment positions, for example, being long and short. Note that our analyses on the measures of sensitivities of industrial sectors and the signs of the factor loadings are conditional on the identification restrictions on the factors and factor loadings.

For the above purposes, all the firms are categorized to one of the ten industrial sectors based on Industry Classification Benchmark (ICB)<sup>17</sup>: (i) *Oil & Gas*; (ii) *Basic Materials*; (iii) *Industrials*; (iv) *Consumer Goods*; (v) *Health Care*; (vi) *Consumer Services*; (vii) *Telecommunications*; (viii) *Utilities*; (ix) *Financials*; (x) *Technology*. Then, for a given factor, the factor loadings are grouped into the negatives and the positives. For each group, the portion of the sum of the absolute value of the factor loadings which belong to each industrial sector is computed and reported. Specifically, we compute the following statistics for factor  $\ell$  and industry  $s$  for given estimation window:

$$T_{b_{\ell},s}^- = \frac{\sum_{i=1}^N \hat{b}_{i\ell} 1\{\hat{b}_{i\ell} < 0\} 1\{i \in s\}}{\sum_{i=1}^N \hat{b}_{i\ell} 1\{\hat{b}_{i\ell} < 0\}}, \quad T_{b_{\ell},s}^+ = \frac{\sum_{i=1}^N \hat{b}_{i\ell} 1\{\hat{b}_{i\ell} > 0\} 1\{i \in s\}}{\sum_{i=1}^N \hat{b}_{i\ell} 1\{\hat{b}_{i\ell} > 0\}}$$

where  $\hat{b}_{i\ell}$  is the estimated factor loading of  $i$ th firm security, and  $1\{A\}$  is the indicator function which takes unity if  $A$  is true and zero otherwise. We regard the portion  $T_{b_{\ell},s}^-$  and  $T_{b_{\ell},s}^+$  as the statistical measure of the negative and positive sensitivities of the  $s$ th industry to the  $\ell$ th factor. The average of the portion of the industrial sectors in S&P500 and the average of  $T_{b_{\ell},s}^-$  and  $T_{b_{\ell},s}^+$  for the four factors over the estimation windows  $\tau = \text{Sept 1998}, \dots, \text{April 2018}$ , are reported in Figure SP2.

Figure SP2(a) shows the portion of the industrial sectors to which the securities consists of S&P500 belong, and the measure  $T_{b_{1},s}^+$  for the first factor. All the loadings to the first factor have the same sign (and it is chosen to be positive), which strongly suggests that this is the market factor. As one might expect, the ‘beta’ (the factor loading) of defensive industries, *Oil&Gas*, *Health Care*, *Telecoms* and *Utilities* is relatively small. The ‘beta’ of cyclical industries such as *Industrials*, *Financials* and *Basic Materials*, is noticeably high. The averages of the measures of negative and positive industrial contributions to the second factor loadings are reported in Figure SP2(b). It shows that *Utility* and *Financials* account for around 43% and 23% of negative loadings, respectively, whilst *Technology*, *Industrials* and *Basic Materials* share 40%, 17% and 14% of positive loadings, respectively. The averages of  $T_{b_{\ell},s}^-$  and  $T_{b_{\ell},s}^+$  for the third factor are reported in Figure SP2(c). It is clear that this is the *Oil&Gas* factor, which share the 67% of the negative loadings. *Financials*, *Consumer Services* and *Consumer Goods* share 29%, 23% and 19% of positive loadings, which means that these industrial sectors move opposite direction to the *Oil&Gas* with respect to the third factor. In view of Figure SP2(d), the dominating industry of the fourth factor is *Utility*, which share 43% of positive loading, together with *Health Care* with 17% of the

<sup>17</sup>Refer to FTSE Russell for more details about ICB.

share. No dominant industry is found for negative loadings, which are equally shared by cyclical industries.

In turn we discuss each factors in more details by analyzing Table SP1, Figures SP1 and SP2. The first factor does seem to be almost always “strong,” in that the absolute sum of factor loadings is proportional to  $N$ . As reported in Table SP1, the average of  $\alpha_1$  over the month windows is 0.995 and standard deviation is very small (0.004) with the minimum value of 0.979. Also as is shown later, all the values of the factor loadings to this factor have the same sign, which strongly suggests that this is the market factor. Now we turn our attention to the rest of the factors. The divergence rates for the rest of the common components,  $\alpha_2$ ,  $\alpha_2$  and  $\alpha_4$ , exhibit very different trajectory over the months, and their orders in terms of value change (i.e., their plots cross).

Let us see the trajectory of  $\alpha_2$ . From Figure SP2(b), under our identification condition, the second factor can be understood of *Utility* and *Financials* versus *Technology*, *Industrials* and *Basic Materials*. In Figure SP1 it is seen that  $\alpha_2$  moves around 0.80 until October 1998, but from this month it sharply goes down and stay below 0.75 to October 1999. Then it sharply goes up to achieve 0.83 in February 2000. Indeed, this period corresponds to the turbulence of *Basic Material* stock index during 1998-2003, the fall of *Industrials* stock index around 2001-2 and the dot com bubble towards the peak in 2000. Since then, during most of the 2000s,  $\alpha_2$  goes above 0.85. After achieving the peak of 0.895 in April 2009, it steadily decreases and stabilizes around 0.75 from November 2012 onward, during which often this factor is not estimated but the fourth factor is.

Now let us analyze the move of  $\alpha_3$ . From Figure SP2(c), under our identification condition, the third factor can be understood of *Oil&Gas* versus *Financials*, *Consumer Services* and *Consumer Goods*. According to Table SP1,  $\alpha_3$  has the lowest average. In Figure SP1, it looks co-moving with  $\alpha_2$ , around 0.1 below, between September 1989 and July 2008. The exceptions are the periods from 1991 to 1992 and from 1999 to 2000, during which  $\alpha_3$  and  $\alpha_2$  are very close. A sharp rise of  $\alpha_3$  is observed from July 2008 to April 2009. This period coincides with the 2008 financial crisis. In just ten months, it goes up by 0.12, from 0.74 to 0.86. This can be interpreted that the *Oil&Gas* industry was sharply affected by the crisis.  $\alpha_3$  exceeds  $\alpha_2$  in December 2010, and this change of the order remains to the latest data point, April 2018.

Now let us analyze the move of  $\alpha_4$ . From Figure SP2(d), under our identification condition, the fourth factor can be understood of *Utility* and *Health Care* versus cyclical industries. As shown in Figure SP1, the first estimate of the fourth factor appears in February 2004, with the value of  $\alpha_4$  being 0.80. Since its appearance, often it is not estimated but it is from March 2010 onward, seemingly becoming more and more stronger toward the latest month, April 2018. Since its first appearance, the value of  $\alpha_4$  is mostly between 0.75 and 0.80. After the sharp one off drop in February 2015,<sup>18</sup>  $\alpha_4$  rises to become the highest next to the first factor from November 2016 onward.

---

<sup>18</sup>This coincides with the period at bottom of the biggest sharp fall in oil price between 2014–2015.

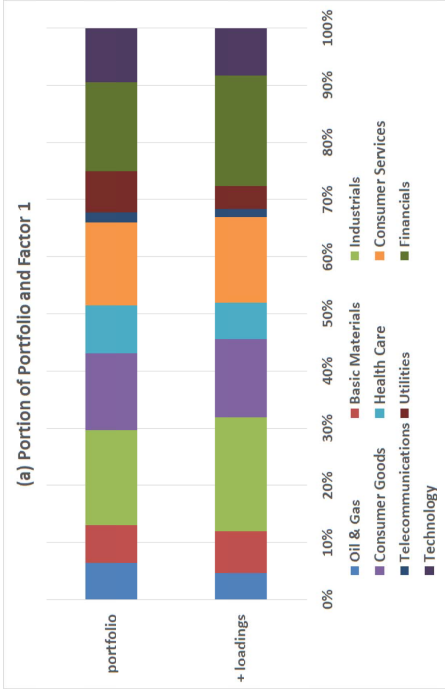


Figure 5: the portion of the industrial sectors in S&P500 and in the Figure 6: the portion of the industrial sectors in the positive/negative second factor loadings

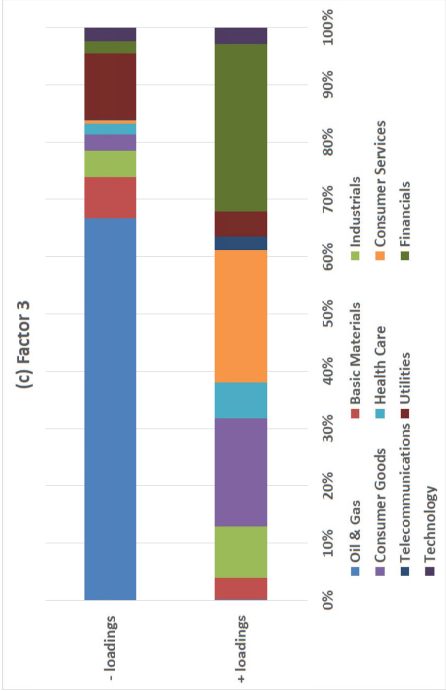
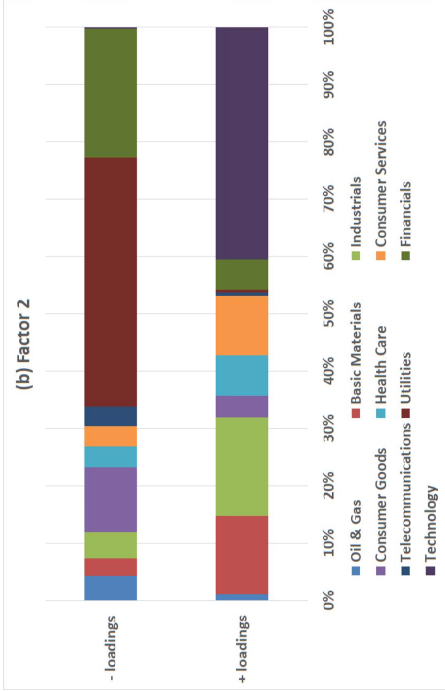


Figure 7: the portion of the industrial sectors in the positive/negative third factor loadings

