

A Generalized Factor Model with Local Factors*

Simon Freyaldenhoven[†]

Brown University

This version: May 31, 2018

Abstract

I extend the theory on factor models by incorporating “local” factors into the model. Local factors only affect an unknown subset of the observed variables. This implies a continuum of eigenvalues of the covariance matrix, as is commonly observed in applications. I derive which factors are pervasive enough to be economically important and which factors are estimable using the common principal component estimator. I then introduce a new class of estimators to determine the number of those relevant factors. Unlike estimators that have been proposed in the past, my estimators are the first to use information in the eigenvectors as well as in the eigenvalues. I find strong evidence for the presence of local factors in a large panel of US macroeconomic indicators.

JEL-Classification: C38, C52, C55

KEYWORDS: high-dimensional data, factor models, weak factors, local factors, sparsity

*I am grateful to Frank Kleibergen, Adam McCloskey, Eric Renault and Jesse Shapiro for their advice. Their insightful comments on earlier versions of this draft have been of great help. I also thank Christian Hansen, Phillip Ketz, Susanne Schennach, audiences at Brown University, the European School of Management and Technology, the Federal Reserve Bank of Philadelphia, the University of Amsterdam, the University of Konstanz, the University of St. Gallen, and the University of Virginia, as well as participants of the 2016 European Conference of the Econometrics Community, the 2016 Maastricht Workshop on Advances in Quantitative Economics, the 2017 CIREQ Econometrics conference and the 2017 NBER/NSF Time Series Conference for helpful comments and suggestions. An earlier draft of this paper was circulated under the title “Factor Models of Arbitrary Strength”.

[†]Department of Economics, Brown University, 64 Waterman Street, Providence, RI 02912.

Email: simon.freyaldenhoven@brown.edu

1 Introduction

Factor models allow for a large number of economic variables to be distilled into a small number of reference variables, enabling the analysis of otherwise prohibitively complex datasets. This paper generalizes standard factor models by introducing a novel theoretical framework incorporating factors of varying strength. Instead of ruling out “local” factors that only affect a subset of the observed variables, as is commonly done in the literature, I explicitly allow for factors of arbitrary strength. Here, the strength of a factor is defined by the number of outcomes it affects.¹ I derive which factors are strong enough to be economically important and which factors are estimable using the common principal component estimator. I then introduce a new class of estimators to determine the number of those relevant factors.

While there exists a multitude of estimators for the number of factors (e.g. Bai and Ng (2002), Onatski (2010), Ahn and Horenstein (2013)), existing estimators are derived from the empirical distribution of the eigenvalues. I argue that in a setting with local factors there is additional information in the eigenvectors and propose to exploit this additional information by incorporating partial sums of the eigenvectors into the estimator.

While local factors have long been acknowledged, the current literature requires a clear distinction between “large” and “small” groups of affected variables, ruling out factors that drive a decreasing fraction of the observables. For example, the handbook chapter of Connor and Korajczyk (1995) distinguishes between factors affecting at most a fixed number of firms and factors affecting at least a constant proportion of all firms.² This paper proposes a more general model that allows for groups of intermediate sizes. This generalization provides a better approximation to the data under a given sample size.

Although the standard model implies a clearly visible separation of the eigenvalues of the covariance matrix into two groups (large eigenvalues representing factor related variation and small eigenvalues representing idiosyncratic variation), such a visible separation is typically not found in practice. For example, a popular dataset in which factor models have been used is the “Stock & Watson” dataset (e.g. Stock and Watson (2002a), De Mol et al. (2008)), consisting of a large panel of US macroeconomic indicators. Figure 1 depicts the distribution of eigenvalues in an updated

¹Note that this is different from the weak factor framework of Onatski (2012) and Kleibergen (2009).

²Specifically, Connor and Korajczyk (1995) state: “Suppose that there is a large number (n) of assets each representing the common shares of one firm. Each firm belongs to one of a large number (m) of industries each with a small number (h , with h approximately equal to n/m) of firms. Idiosyncratic returns are correlated within industries but uncorrelated across industries. [...] Holding h constant and letting n and m increase, this series of covariance matrices has bounded eigenvalues. [...] On the other hand, suppose that there is a small number, k , of sectors, each containing n/k firms. All firms within sector j are subject to sector shock f_j with unit betas (for simplicity). Firms in sector j are unaffected by the shocks of other sectors. Given these assumptions, the sector shocks constitute pervasive risk. Note the clear distinction between industries (a small proportion of the firms are in each industry) versus sectors (a substantial proportion of the firms are in each sector).”

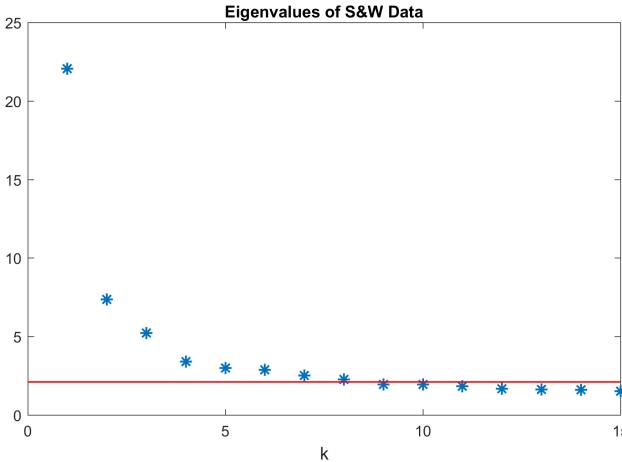


Figure 1: 20 largest eigenvalues of the covariance matrix for a dataset of 94 macroeconomic indicators in the US. Solid line indicates cutoff chosen according to Bai and Ng (2002) (with $r_{max} = 15$) to determine the number of factors. A more detailed discussion of this application can be found in section 6.

vintage of this data.³ The solid line indicates the cutoff between the two groups as chosen by a criterion of Bai and Ng (2002). A model implying a continuum of eigenvalues represents the data much better than such a classification into two groups. In finance, empirical studies on Arbitrage Pricing Theory (Ross (1976)) similarly point to a continuum of factor strengths. For example, in a cross section of asset returns Trzcinka (1986) finds that, while the first eigenvalue dominates, the first 6 eigenvalues diverge at differing rates.

There are two kinds of weak factors that may induce a continuous decay in the distribution of eigenvalues. Either a weaker factor can have a weak effect on all observables, or it can affect only a subset of observables. This paper will largely focus on the latter scenario, although some results extend to the former. Recently a number of empirical studies have also postulated a structure with group-specific factors, sometimes called hierarchical factor models (e.g. Boivin and Ng (2006), Moench et al. (2013), Dias et al. (2013)). The findings in De Mol et al. (2008) further provide empirical evidence for such a structure (see the discussion in Freyaldenhoven (2017)). Ando and Bai (2017) also considers group-specific factors, but requires all group sizes to be comparable to the overall cross-sectional dimension. To the best of my knowledge, the only theoretical paper in the direction of local factors in the sense of this paper is Wang (2008). However, unlike Wang (2008), I do not require the group structure and factors strengths to be known to the practitioner *a priori*.

Onatski (2009), Onatski (2010) and Onatski (2012) propose a framework for weak factors through random matrix theory, and a similar model to the one used in this paper has been con-

³I discuss the data in more detail in section 6.

sidered in the large body of literature on sparse PCA (e.g. Paul and Johnstone (2012), Cai et al. (2013)). These papers typically build on stronger assumptions on the error terms, assume bounded eigenvalues of the covariance matrix, and remain largely agnostic about the factors themselves. By considering a diverging eigenvalue regime, and explicitly modeling the factors, we are able to impose less restrictive assumptions on the error structure. There is also a related literature on sparse factor models under a Bayesian framework (e.g. Carvalho et al. (2008), Gao et al. (2013) and Pati et al. (2014)). Finally, by considering a continuum of factor strengths, this paper is similar in spirit to the extensive literature in econometrics on identification with varying convergence rates (e.g. Andrews and Cheng (2012), Antoine and Renault (2012)).

Before I formally introduce the model, the following are some concrete examples of economic models to which this paper applies.

Example 1. Arbitrage Pricing Theory.

Consider an unobserved common shock that affects only a subset of the population at the company level, for example, a new law that affects only large firms. As the number of firms, n , increases, one reasonable assumption is that the number of large firms increases at a rate slower than n (Chudik et al. (2011)). Unlike traditional factor models, the framework of this paper allows for this and is in line with the empirical finding indicating that the largest eigenvalues of the sample covariance matrix of asset returns diverge at differing rates (e.g. Trzcinka (1986)) .

Empirical evidence on whether weaker factors are priced appears somewhat mixed (e.g. Shukla and Trzcinka (1990)). In section 3.2 I use the results in Green and Hollifield (1992) adapted to my framework to derive theoretical bounds on the strength of factors that will be priced. I find that the number of factors which are priced depends directly on the degree of diversification of the portfolios on the efficient frontier. The better diversified these portfolios are, the smaller the number of factors that have a non zero factor premium.

Example 2. The origins of aggregate fluctuations.

There is an ongoing debate about the origins of fluctuations in the aggregate economy (see, e.g. Foerster et al. (2011)). Long and Plosser (1983) suggest that sectoral shocks may account for GDP fluctuations. With a fixed number of sectors, these sectoral shocks affect a fixed proportion of firms and can be viewed as aggregate shocks themselves. In contrast, Horvath (1998) investigates conditions under which an economy with n sectors can have a volatility that does not decay according to $\frac{1}{\sqrt{n}}$. By modeling sectoral shocks as local factors affecting the corresponding subset of firms, this can be mapped into the framework of this paper. I show in Section 3.2 that in an economy with n firms \sqrt{n} -convergence for the aggregate growth rate of the economy fails when there are sectoral shocks affecting proportionally more than \sqrt{n} firms. I therefore find that aggregate fluctuations can be attributed to sectors proportionally larger than \sqrt{n} firms.

Example 3. *Macroeconomic forecasting.*

In a widely cited paper Boivin and Ng (2006) investigate the properties of the principal components estimators in finite samples. Specifically, they document conditions under which adding more data can be undesirable for factor estimation.

As a stylized model they consider macroeconomic panels with two factors. Some series are driven by two factors, some are only affected by one factor, while others are not associated with any factor. For example, the first n_1 series (only affected by the first factor) might be output and employment type series, the next n_2 series might be prices (affected by F_2), the following n_3 series represent interest rates and are affected by both, while variations in the remaining series are purely idiosyncratic. If, for example, the cross section contains relatively few series representing prices and interest rates, this fits the framework of this paper.

Boivin and Ng (2006) use a Monte Carlo study to establish that the performance of the principal component estimator deteriorates as more “noisy” series are added, effectively making factors local in the sense of this paper. I provide an analytical framework, tying the convergence rate of a factor estimate to the factor’s strength, that can help to explain their result.

2 A Model with Local Factors

To set up notation, define an n -dimensional process by X_t , $t = 1, 2, \dots, T$. Let F_k , $k = 1, 2, \dots, r$ denote the true factors. $\Lambda = [\lambda_{\cdot 1} \lambda_{\cdot 2} \cdots \lambda_{\cdot r}] = [\lambda_1 \lambda_2 \cdots \lambda_n]'$ denotes the matrix of factor loadings. Throughout, I use the running indices s and t for the T observations, indices i, j for the n variables, and k and l for the r factors. I assume that the data has a static factor structure⁴:

$$\underset{(T \times n)}{X^{(n)}} = \underset{(T \times r)(r \times n)}{F^{(n)}} \underset{(r \times n)}{\Lambda^{(n)'}} + \underset{(T \times n)}{e^{(n)}}. \quad (1)$$

I treat both Λ and F as parameters of the distribution of X . In cases when factors are random, such an approach is equivalent to conditioning on a particular realization of F . All results should still hold if F and Λ are random, provided they are independent of each other and both $E\|F_t\|^4$ and $E\|\lambda_i\|^4$ are bounded.

I denote the p th largest eigenvalue of a matrix A by $\psi_p(A)$ and the Frobenius norm of a matrix B by $\|B\|$, such that $\|B\|^2 = \text{tr}(B'B) = \sum_{ij} b_{ij}^2$. I further make extensive use of the notion that certain quantities diverge at particular rates and write $a_n \asymp b_n$ for two sequences a_n, b_n if $a_n = O(b_n)$ and $b_n = O(a_n)$. I write $Y_n = \bar{O}_p(n^\gamma)$ as short hand notation for $Y_n = O_p(\min\{1, n^\gamma\})$.

⁴A more general setup would be the dynamic factor model of Forni et al. (2000), allowing for factor loadings that are represented by possibly infinite lag polynomials. However, whenever the order of such lag polynomials is bounded, the model can be rewritten in the static form with constant factor loadings, where the factors are augmented by a set of their own lags (See, e.g. Stock and Watson (2006))

Finally let ι_p denote a vector with a 1 at entry p and zeros everywhere else, with the dimension varying, but obvious from context.

Assumption 1. *There exist positive constants c, C and a diagonal matrix $D_r^{(n)}$ with diagonal entries $d_1^{(n)}, d_2^{(n)}, \dots, d_r^{(n)}$, such that as $n, T \rightarrow \infty$:*

$$(a) \ n/T \rightarrow c$$

$$(b) \ \Lambda^{(n)'} \Lambda^{(n)} = D_r^{(n)}, \ d_1^{(n)} > d_2^{(n)} > \dots > d_r^{(n)} \text{ and } |\lambda_{ik}| < C \ \forall i$$

$$(c) \ \frac{1}{T} F^{(n)'} F^{(n)} = I_r \text{ and } |F_{tk}| < C \ \forall t$$

Part (a) of Assumption 1 requires n and T to be comparable even asymptotically. This assumption is common in the literature (e.g. Onatski (2012), Ahn and Horenstein (2013)) and plausible in many applications of factor models. Part (b) of Assumption 1 is considerably weaker than the standard assumptions in the literature (e.g. Stock and Watson (2002a), Bai and Ng (2002), Bai (2003), Ahn and Horenstein (2013)) in that the entries in $D^{(n)}$ are not assumed to diverge proportionally to n , thus allowing for weaker factors. All entries in $D^{(n)}$ can have different rates. Thus, one can think of Assumptions 1(b)-(c) as simply normalizations or identifying restrictions.

To simplify notation, I will omit the superscript (n) on matrices X, Λ, F, D and e in what follows.

Assumption 2. *For each factor k , the entire set of indices $i = 1, 2, \dots, n$ can be partitioned into a set of indices \mathcal{A}_k with cardinality $|\mathcal{A}_k| \asymp n^{\alpha_k}$ for some $\alpha_k \in [0, 1]$ and its complement such that, as $n, T \rightarrow \infty$ for all k :*

$$(a) \ \sum_{i \in \mathcal{A}_k} \lambda_{ik}^2 \asymp n^{\alpha_k}$$

$$(b) \ \sum_{i \notin \mathcal{A}_k} \lambda_{ik}^2 < C \text{ for some } C < \infty$$

Assumption 2 allows for the loadings of any given factor k to be concentrated on an asymptotically vanishing fraction of variables. It states that any given factor fulfills the conventional pervasiveness assumption only on an unknown subset of all outcomes (\mathcal{A}_k), while the remaining loadings are small in the sense that their squares are summable. As a specific example, consider a cross section of n assets and an industry with a size proportional to \sqrt{n} of the assets. Suppose there exists an industry-specific factor F_l that affects only those assets. Then, $\alpha_l = .5$ and $\sum_{i \notin \mathcal{A}_l} \lambda_{il}^2 = 0$, such that Assumption 2 holds. The standard assumptions in the literature correspond to assuming $\alpha_k = 1$ for all factors, thus ruling out any such local factors.

Assumption 3. *There exist constants $c > 0, C < \infty$ and a constant $d \in (0, 1]$ (which may depend on c), such that*

- (a) $E(e_{ti}) = 0, E|e_{ti}|^4 \leq C$
- (b) $\sum_{t=1}^T |E\left(\frac{e'_s e_t}{n}\right)| \leq C \forall s \quad \text{and} \quad \sum_{j=1}^n |E\left(\frac{e'_i e_j}{T}\right)| \leq C \forall i$
- (c) for every (t,s) , $E|\frac{1}{\sqrt{n}}[e'_s e_t - E(e'_s e_t)]|^4 \leq C$
- (d) $E\|\frac{1}{\sqrt{nT}}\sum_{s=1}^T F_s[e'_s e_t - E(e'_s e_t)]\|^2 \leq C \forall t$
- (e) $\psi_1\left(\frac{e'_e}{T}\right) = O_p(1)$ and $P\left(\psi_{[dn]}\left(\frac{e'_e}{T}\right) \geq c\right) = 1$ for some $d > 0$

Assumption 4. For any $k, l < r$:

$$(a) \frac{\Lambda'_{k} e_t}{n^{\frac{1}{2} \alpha_k}} = O_p(1) \quad \forall t$$

$$(b) \frac{\Lambda'_{k} e' F_{l}}{n^{\frac{1}{2} \alpha_k} T^{\frac{1}{2}}} = O_p(1)$$

Assumptions 3 and 4 concern the possibly correlated noise. Assumption 3 rules out that there is too much dependence in the error terms and is standard in the literature (Bai (2003), Bai and Ng (2006b)). More primitive conditions can be provided that imply part (e) (see Onatski (2015), Moon and Weidner (2017)). Assumption 4 is weaker than one that requires a number of Central Limit Theorems to hold. With $\alpha_k = 1$ for $k = 1, \dots, r$, it is implied by Assumptions F2 and F3 in Bai (2003).

Remark 1. Let $r_1 + r_2 = r$, $\alpha_k > \tau$ for $k = 1, \dots, r_1$ and $\alpha_k \leq \tau$ for $k = r_1 + 1, \dots, r$ for some fixed value of $\tau \in [0, 1]$. In words: Let r_1 be the number of factors affecting proportionally more than n^τ variables, while the remaining factors are less pervasive. We can then rewrite the factor structure (1) as

$$\underset{(T \times n)}{X} = \underset{(T \times r)(r \times n)}{F} \underset{(T \times n)}{\Lambda'} + \underset{(T \times n)}{e} \tag{2}$$

$$= \underset{(T \times r_1)(r_1 \times n)}{F^s} \underset{(T \times n)}{\Lambda^{s'}} + \underset{(T \times r_2)(r_2 \times n)}{F^w} \underset{(T \times n)}{\Lambda^{w'}} + \underset{(T \times n)}{e} \tag{3}$$

$$= \underset{(T \times r_1)(r_1 \times n)}{F^s} \underset{(T \times n)}{\Lambda^{s'}} + \underset{(T \times n)}{u}, \tag{4}$$

such that the weakest r_2 factors are incorporated into the error term. If we think of this as a model with effectively r_1 factors, it follows that $\psi_1(uu'/n)$ is no longer bounded. We can therefore think of Assumptions 1-4 as a generalization of standard factor models in two ways: they allow for the presence of weaker factors than were previously allowed for in the literature and they allow for a stronger dependence in the errors. By including more (weaker) factors, a practitioner can choose how much of the correlation among the observables she wishes to explicitly model. In fact one can

generally always include additional factors, even if the corresponding eigenvalue will be bounded. One can therefore think of r as an upper bound on the number of factors in this paper.

However, a framework of increasingly weaker (more local) factors and the resulting continuum of factor strengths also immediately raises the question of how many factors a practitioner should keep in the model. We model this choice of r_1 through the complexity parameter τ . A practitioner chooses a threshold $\tau \in [0, 1]$ to indicate a lower bound on the strength of the factors she wishes to keep in the model. While my methods allow for any choice of τ , I discuss this choice in sections 3.1-3.2.

Although I treat r as fixed, thus not allowing the number of factors to grow with the sample size, conceptually, my framework would allow for this. In contrast, with all factors affecting a non-vanishing fraction of the observables, as is the case in standard factor models, a growing number of factors is generally impossible⁵.

Finally, I denote $\alpha_{max} = \alpha_1$ and note that all auxiliary lemmata for the proofs in the following sections are relegated to the appendix.

3 Weak Asymptotics

I first show what the introduction of local factors implies for the empirical distribution of the eigenvalues of the matrix $\frac{X'X}{T}$. This is the quantity depicted in Figure 1 and often included in applications to justify the use of a factor model. I start with the following lemma:

Lemma 1. *Under Assumptions 1 and 2:*

$$\psi_k\left(\frac{\Lambda F' F \Lambda'}{T}\right) \begin{cases} \asymp n^{\alpha_k}, & k = 1, 2, \dots, r \\ = 0 & k > r. \end{cases} \quad (5)$$

Proof. If $k \leq r$:

$$\psi_k\left(\frac{\Lambda F' F \Lambda'}{T}\right) = \psi_k(\Lambda \Lambda') = \psi_k(\Lambda' \Lambda) \quad (6)$$

$$= \sum_{i=1}^n \lambda_{ik}^2 = \sum_{i \in \mathcal{A}_k} \lambda_{ik}^2 + \sum_{i \notin \mathcal{A}_k} \lambda_{ik}^2 \asymp n^{\alpha_k} + O_p(1) \quad (7)$$

$$\asymp n^{\alpha_k}, \quad (8)$$

⁵With X standardized and no correlation in the error terms, $\text{Corr}(X_i, X_j) = \sum_{k=1}^r \lambda_{ik} \lambda_{jk}$. Thus, if r is an increasing sequence, the correlation between any two observables would continue to increase as more factors are added, unless we model the loadings as drifting towards zero. Allowing the number of factors to grow with the sample size is left as an interesting extension for future research.

where the equality in the second line follows from Assumption 2.

If $k > r$: the result immediately follows from the fact that $\text{rank}(\Lambda F' F \Lambda') = r$. \square

The properties of the eigenvalues of the matrix $\frac{X'X}{T}$ then follow:

Theorem 1. *For any given Factor k ($k = 1, 2, \dots, r$), under Assumptions 1-3:*

$$\psi_k\left(\frac{XX'}{T}\right) \begin{cases} \asymp n^{\alpha_k} \text{ for } k = 1, 2, \dots, r \\ = O_p(1) \text{ for } k = r + 1, \dots, n. \end{cases} \quad (9)$$

Proof. By the singular value version of Weyl's inequalities (Horn and Johnson (2012), p.454):

$$\sigma_{k+l-1}(A + B) \leq \sigma_k(A) + \sigma_l(B) \quad 1 \leq k, l \leq q, \quad k + l \leq q + 1, \quad (10)$$

where $\sigma_k(A)$ denotes the k 'th largest singular value of a matrix A . Therefore, with $A = F\Lambda'$, $B = e$ and $l = 1$, for $k = 1, 2, \dots, r_{\max}$:

$$\sigma_k(X) \leq \sigma_k(F\Lambda') + \sigma_1(e). \quad (11)$$

Since $\sigma_k(A) = \sqrt{\psi_k(AA')}$ for any matrix A , it follows that

$$\sqrt{\psi_k(XX')} \leq \sqrt{\psi_k(F\Lambda'\Lambda F')} + \sqrt{\psi_1(ee')}. \quad (12)$$

And I therefore conclude, using Lemma 1 and Assumption 3(e) respectively for the two eigenvalues on the RHS:

$$\psi_k\left(\frac{XX'}{T}\right) \leq \psi_k\left(\frac{F\Lambda'\Lambda F'}{T}\right) + \psi_1\left(\frac{ee'}{T}\right) + 2\sqrt{\psi_k\left(\frac{F\Lambda'\Lambda F'}{T}\right)}\sqrt{\psi_1\left(\frac{ee'}{T}\right)} \quad (13)$$

$$\leq C_1 n^{\alpha_k} + O_p(1) + O_p(n^{\frac{1}{2}\alpha_k}) \quad (14)$$

$$\leq C_2 n^{\alpha_k}. \quad (15)$$

Similarly, again by Weyl's inequalities:

$$\sigma_k(X - e) \leq \sigma_k(X) + \sigma_1(-e) \quad (16)$$

$$\sigma_k(F\Lambda') \leq \sigma_k(X) + \sigma_1(e) \quad (17)$$

$$\sqrt{\psi_k\left(\frac{F\Lambda'\Lambda F'}{T}\right)} \leq \sqrt{\psi_k\left(\frac{XX'}{T}\right)} + \sqrt{\psi_1\left(\frac{ee'}{T}\right)} \quad (18)$$

$$\sqrt{\psi_k\left(\frac{XX'}{T}\right)} \geq \sqrt{c_1 n^{\alpha_k}} - O_p(1) \quad (19)$$

and I therefore also conclude that $\psi_k(\frac{XX'}{T}) \geq c_2 n^{\alpha_k}$. \square

Under a scenario with r strong factors ($\alpha_k = 1$ for all $k = 1, 2, \dots, r$) this reduces to the standard result in the literature: the first r eigenvalues diverge at rate n (Connor and Korajczyk (1993), Bai and Ng (2002), Hallin and Liska (2007)). I extend this result to allow for weaker factors with the slower divergence rates of Theorem 1 for factors that affect only a subset of the observed variables. Note that we can replace Assumption 2 with the high level assumption $\sum_{i=1}^n \lambda_{ik}^2 \asymp n^{\alpha_k}$ and Theorem 1 still holds. The result in Theorem 1 therefore extends to weak factors in general and does not need the sparsity pattern that is imposed by Assumption 2.

Theorem 1 provides a possible explanation for the continuum of eigenvalues often observed, as in Figure 1. While conventional factor models imply a large gap in the eigenvalue distribution after the r th eigenvalue, the eigenvalues corresponding to local factors will fall into this gap.

Recall the earlier distinction of factors into two groups: $F = [F_1, \dots, F_{r_1}, F_{r_1+1}, \dots, F_r] = [F^s, F^w]$, such that $r = r_1 + r_2$, $\alpha_k > \tau$ for $k = 1, 2, \dots, r_1$ and $\alpha_k \leq \tau$ for $k = r_1 + 1, \dots, r$ for some user specified threshold $\tau \in [0, 1]$. To provide guidance on how to choose the tuning parameter τ (the lower bound on the pervasiveness of factors one wishes to keep in the model), I next consider the following two questions:

1. When is a factor strong enough to be estimated consistently?
2. When is a factor strong enough to be of interest in some common economic models?

3.1 The Principal Component Estimator

I will begin with the first question and consider the standard estimator in the literature: estimation of both the factors and their loadings is achieved through the principal component estimator (see Stock and Watson (2002a), Bai and Ng (2002), Bai (2003)). I obtain the following theorem:

Theorem 2. *Let \hat{F}_k be defined as the standardized eigenvector corresponding to the k th largest eigenvalue of $\frac{XX'}{n}$. Then, under assumptions 1-4,*

$$\hat{F}_{tk} - F_{tk} = O_p(n^{1-2\alpha_k}) + O_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}). \quad (20)$$

Proof. First define a matrix H as follows:

$$H = \Lambda' \Lambda \frac{F' \hat{F}}{T} \hat{D}_K^{-1}, \quad (21)$$

where \hat{D}_K is a diagonal matrix with the K largest eigenvalues of $\frac{X'X}{T}$ on the main diagonal. By Lemma 9: $H_{.k} = \iota_k + O_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2}-\alpha_k})$.

Consequently, combining this with Lemma 8:

$$\hat{F}_{kt} - F_{kt} = (\hat{F}_{kt} - H'_{k\cdot} F_t) + (H'_{k\cdot} - \iota'_k) F_t \quad (22)$$

$$= O_p(n^{1-2\alpha_k}) + O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) + O_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2}-\alpha_k}) \quad (23)$$

$$= O_p(n^{1-2\alpha_k}) + O_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2}-\alpha_k}). \quad (24)$$

□

I note that this result may be of interest to a practitioner for two reasons. First, this establishes a lower bound in terms of factor strength for consistency of the principal component estimator ($\alpha_k > \frac{1}{2}$). Further, even for factors that are estimated consistently, it suggests that the estimation of a factor becomes worse as its strength decreases. (This is documented in simulations in Boivin and Ng (2006).) The intuition is clear: As fewer cross-sections carry a signal about F_k , the precision of its estimate decreases. However, the fact that weaker factors tend to be estimated with less precision seems to be largely unaccounted for in the current literature⁶. In cases in which factor estimates are used that correspond to weaker factors, Theorem 2 at least suggests to be cautious with respect to the standard errors of these estimates.

I also obtain a similar result for the factor loadings:

Theorem 3. Let $\hat{\Lambda}' = \frac{\hat{F}'X}{T}$, with \hat{F} defined as before. Then, under assumptions 1-4:

$$\hat{\lambda}_{ik} - \lambda_{ik} = \bar{O}_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}). \quad (25)$$

Proof.

$$\hat{\lambda}_{ik} = \frac{\hat{F}'_k X_i}{T} = \frac{1}{T} \hat{F}'_k F \lambda_i + \frac{1}{T} \hat{F}'_k e_i \quad (26)$$

$$= \lambda_{ik} + \left(\frac{\hat{F}'_k F}{T} - \iota'_k \right) \lambda_i + \frac{1}{T} (\hat{F}_k - F_k)' e_i + \frac{1}{T} F' e_i \quad (27)$$

$$= \lambda_{ik} + \left(\frac{\hat{F}'_k F}{T} - \iota'_k \right) \lambda_i + \frac{1}{T} (\hat{F}_k - F H_{\cdot k})' e_i + \frac{1}{T} (F H_{\cdot k} - F_k)' e_i + \frac{1}{T} F' e_i \quad (28)$$

$$= \lambda_{ik} + \left(\frac{\hat{F}'_k F}{T} - \iota'_k \right) \lambda_i + \frac{1}{T} (\hat{F}_k - F H_{\cdot k})' e_i + \frac{1}{T} (H_{\cdot k} - \iota_k)' F' e_i + \frac{1}{T} F' e_i \quad (29)$$

$$= \lambda_{ik} + O_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + O_p(n^{1-2\alpha_k}), \quad (30)$$

where the last equality follows from Lemmata 7, 9 and 12 as well as assumption 4(b). Since also $\frac{\hat{F}'_k X_i}{T} \leq \frac{1}{T} \|\hat{F}\| \|X_i\| = O_p(1)$, this completes the proof. □

⁶The exception is Onatski (2012), but his setup is quite distinct from the framework in this paper.

n	81	256	625	1296
$ \mathcal{A}_2 = n^{1/4}$	0.1115	0.0772	0.0694	0.0639
$ \mathcal{A}_2 = n^{1/2}$	0.492	0.4848	0.406	0.2872
$ \mathcal{A}_2 = n^{3/4}$	0.7452	0.8528	0.9094	.9392

Table 1: Correlation between estimate and true factor for differing factor strengths as sample size increases. Table based on 400 repetitions. Smaller models not nested in larger ones. See Appendix for exact description of DGP.

Thus, I obtain convergence of the principal component estimator for both the factors and the loadings as long as $\alpha_k > .5$. The following table provides an empirical test to assess the adequacy of these asymptotic results in approximating finite sample properties. Two factors were created, one strong and one weak. The strength of the weak factor is varied from $n^{.25}$ to $n^{.75}$ and the correlation of the estimate with its true counterpart is depicted in Table 1. This correlation can be thought of as a measure of consistency of the PC estimator for the k th factor (Bai (2003), Onatski (2012)). The exact DGP can be found in the Appendix. The numbers in Table 1 are in line with my theoretical findings: Theorem 2 suggests that the correlation between F_k and \hat{F}_k will approach 1 when $\alpha_k > .5$. This corresponds to the third row in Table 1. Additionally, I observe that the estimation error seems to indeed dominate the signal when the factor strength is below this threshold.

In many cases factor models are only estimated as a first step. For example, once they are estimated, the extracted factors are often used as what are usually called diffusion indices in forecast models (see, e.g. (Stock and Watson, 2002a), Bai and Ng (2006a)). From my discussion in the previous section, a natural choice for the number of factors to be included will then be those affecting at least $K_n \asymp \sqrt{n}$, as this marks the threshold for which we established consistency of the PC estimator.

3.2 Which Factors Matter?

Aside from practical issues in estimating factors that only affect a small subset of the cross sections, it is also not clear which factors are of interest to a practitioner in the first place. The following section aims to shed some light on this issue.

To this end I next present two theoretical arguments why factors affecting proportionally more than \sqrt{n} cross sections ($\tau = .5$) will be the natural target in many instances. One is derived from the Arbitrage Pricing Theory of Ross (1976) and a second argument builds on the works of Gabaix (2011) in the context of aggregate fluctuations in the economy. Note that for the two examples that follow I treat the Factors as random in line with the literature in those fields.

3.2.1 Arbitrage Pricing Theory

I assume that the n -vector of demeaned asset returns $R_t - \mathbb{E}(R_t)$ for a given t follows a factor structure with potentially local factors as in as in the previous sections

$$R_i - \mathbb{E}(R_i) = \underset{(1 \times r)(r \times 1)}{\lambda'_i} F + \underset{(1 \times 1)}{e_i} = \underset{(1 \times K)(K \times 1)}{\lambda_i^{K'}} F^K + e_i^K, \quad (31)$$

treating the the factors as random and the errors as uncorrelated with the factors. Equation (31) emphasizes that in the more general framework of this paper, we can always choose to move some of the weaker factors into the error structure at the expense of more correlation in the error term. Denote the return of a portfolio by $R^p = \sum_i^n w_i R_i$, with $\sum_i^n w_i = 1$. I formalize the term “well-diversified” by imposing a bound on the sup-norm of the weights: $|w_i| \leq W_n \forall i$ ⁷. Following Green and Hollifield (1992), I say that exact APT pricing holds if the mean returns are in the span of the factor loadings and a constant vector:

$$E(R_j) = (1 - \sum_k^K \lambda_{jk}) E(R_0^*) + \sum_k^K \lambda_{jk} E(R_k^*), \quad (32)$$

where the portfolios $R_k^*, k = 0, \dots, K^*$ are “factor-mimicking” portfolios. Their construction is detailed in the appendix and conditions for their existence are given in Huberman et al. (1987). Similarly, I define exact APT to hold in the limit, if, as n increases, there exist sequences of feasible factor-mimicking portfolios R_{nk}^* , such that for any fixed j

$$\lim_{n \rightarrow \infty} E(R_j) - [(1 - \sum_k^K \lambda_{jk}) E(R_{n0}^*) + \sum_k^K \lambda_{jk} E(R_{nk}^*)] = 0. \quad (33)$$

Finally denote by ν_n the return on the global minimum variance portfolio when there are n assets and assume that the mean-variance frontier does not become vertical in the limit, such that there remains a meaningful trade-off between mean and variance.⁸ I then obtain the following proposition:

Proposition 1. *Consider the sequence of efficient (minimum variance) portfolios for some mean return $\mu \neq \lim_{n \rightarrow \infty} \nu_n$. If*

- (i) $W_n = o(\frac{1}{n^\gamma})$, $\tau > \frac{1}{2}$ for every such portfolio, and

⁷While Chamberlain (1983) defines portfolio diversification through the ℓ_2 -norm, the norm that proves tractable here is the sup-norm. Either definition formalizes the idea that the weights on individual assets get small as the universe of assets expands.

⁸This is the equivalent of the “absence of arbitrage” assumption in the Hilbert space setting of Chamberlain and Rothschild (1983).

$$(ii) \lim_{n \rightarrow \infty} \sum_{i=1}^n |Cov(e_i, e_j)| = O(\sqrt{n}),$$

then exact APT pricing holds in the limit with respect to the strongest K factors, where K is defined such that $\alpha_k > \gamma$ for $k = 1, 2, \dots, K$ and $\alpha_k \leq \gamma$ for $k \geq K + 1$.

The proof can be found in appendix D.3 and largely follows the proof of Theorem 3 in Green and Hollifield (1992).

Proposition 1 states that exact APT holds in the limit if the efficient portfolios are well diversified. Further, the number of factors that are priced depends directly on the degree of diversification of the portfolios on the efficient frontier. The better diversified these portfolios are, the smaller the number of factors that have a non zero factor premium.

In particular, with $W_n = o(\frac{1}{\sqrt{n}})$, which yields diversification in the sense of Chamberlain and Rothschild (1983) and Chamberlain (1983), Proposition 1 establishes that exact APT pricing holds in the limit with respect to the r_1 factors affecting proportionally more than \sqrt{n} of the assets (factors with $\alpha_k > .5$).

Proposition 1 holds under more general conditions than the approximate factor model of Chamberlain and Rothschild (1983): I do not require all eigenvalues of the error covariance matrix to be bounded, but explicitly allow for additional, weaker factors. Instead of ruling out the existence of such weaker factors, Proposition 1 establishes that they will not be priced.

3.2.2 Aggregate Fluctuations in the Economy

Consider a simple “Islands” economy with n firms as in Gabaix (2011). Firm i produces a quantity S_{it} of the consumption good. Instead of modeling firm level growth rates as unrelated, I model them as a combination of r mutually independent shocks that may affect several firms, on top of the idiosyncratic shocks. Firm i thus experiences a growth rate equal to

$$\frac{\Delta S_{i,t+1}}{S_{it}} = \frac{S_{i,t+1} - S_{it}}{S_{it}} = \lambda_i F_{t+1} + \sigma_i \varepsilon_{i,t+1}, \quad (34)$$

where σ_i is firm i ’s volatility, and the $\varepsilon_{i,t+1}$ are uncorrelated random variables with mean zero and variance 1. Firms’ growth rates may be correlated through the presence of the first component. However, I do not impose the factors to be pervasive and likely $\lambda_{ik} = 0$ for most firm-factor combinations. Intuitively, these factors can correspond to economy wide shocks, but also sector shocks or the introduction of policies only affecting a subset of firms including shocks that affect as few as two firms. Thus (34) is quite general.

In this stylized model GDP growth is given by:

$$\frac{\Delta Y_{t+1}}{Y_t} = \frac{1}{Y_t} \sum_{i=1}^n \Delta S_{i,t+1} = \sum_{i=1}^n \frac{S_{it}}{Y_t} [\lambda_i F_{t+1} + \varepsilon_{i,t+1}] \quad (35)$$

$$= \sum_{i=1}^n \frac{S_{it}}{Y_t} \lambda_i F_{t+1} + \sum_{i=1}^n \frac{S_{it}}{Y_t} \varepsilon_{i,t+1}. \quad (36)$$

It follows that the variance of GDP growth at time $(t+1)$ conditional on time t information is equal to

$$Var_t \left(\sum_{i=1}^n \frac{S_{it}}{Y_t} \lambda_i F_{t+1} + \sum_{i=1}^n \frac{S_{it}}{Y_t} \varepsilon_{i,t+1} \right) = Var_t \left(\sum_{i=1}^n \frac{S_{it}}{Y_t} \lambda_i F_{t+1} \right) + Var_t \left(\sum_{i=1}^n \frac{S_{it}}{Y_t} \varepsilon_{i,t+1} \right) \quad (37)$$

$$= Var_t \left(\sum_{i=1}^n \frac{S_{it}}{Y_t} \sum_{k=1}^r \lambda_{ik} F_{k,t+1} \right) + \sum_{i=1}^n \left(\frac{S_{it}}{Y_t} \right)^2 \sigma_i^2. \quad (38)$$

For ease of notation, consider firms of equal size ($S_{it} = \frac{Y_t}{n}$) and identical standard deviation ($\sigma_i = \sigma$), and normalize the factors such that $Var(F_{kt}) = 1$. Further assume that, for a given k , the factor loadings are 1 on a subset of size $|\mathcal{A}_k| \asymp n^{\alpha_k}$ and zero everywhere else⁹. Then:

$$Var_t \left(\frac{\Delta Y_{t+1}}{Y_t} \right) = \sum_{k=1}^r \left(\sum_{i \in \mathcal{A}_k} \frac{1}{n} \right)^2 + \sum_{i=1}^n \frac{1}{n^2} \sigma^2 \quad (39)$$

$$\asymp \sum_{k=1}^r n^{2\alpha_k - 2} + \frac{\sigma^2}{n}. \quad (40)$$

It immediately follows that, absent any factors ($r = 0$), $\sigma_{GDP} = \sqrt{Var_t(\frac{\Delta Y_{t+1}}{Y_t})} = \frac{\sigma}{\sqrt{n}}$, which is the reason macroeconomists often appeal to aggregate shocks, since idiosyncratic fluctuations disappear in the aggregate at rate \sqrt{n} . Next consider an economy with r shocks, where r_1 is

⁹Defining the loadings instead in a more general way as in Assumption 2 does not alter any conclusions

the number of factors with $\alpha_k > .5$:

$$\begin{aligned} Var_t \left(\frac{\Delta Y_{t+1}}{Y_t} \right) &\asymp \sum_{k=1}^{r_1} n^{2\alpha_k - 2} + \sum_{k=r_1+1}^r n^{2\alpha_k - 2} + \frac{\sigma}{n} \\ &= \sum_{k=1}^{r_1} n^{2\alpha_k - 2} + O_p\left(\frac{1}{n}\right) \end{aligned}$$

Equation (40) establishes that the important shocks are those with $\alpha_k > \frac{1}{2}$ and that the standard rate of convergence breaks down whenever shocks exist that affect more than \sqrt{n} firms.

This is in line with the granularity conditions derived in Gabaix (2011), who considers heterogeneous firm sizes that may grow with n . Intuitively, with the growth rate of the economy given by the sum of both the idiosyncratic and factor shocks in my context, we can think of the sector shocks as additional, but larger firms. Then the economy consists of $n + r$ components (with $r \ll n$). Proposition 2 in Gabaix (2011) establishes that $\sigma_{GDP} \asymp \frac{1}{\sqrt{n}}$ only if the largest firm has a relative weight of at most $W_n = O\left(\frac{1}{\sqrt{n}}\right)$. Of course this corresponds exactly to the limit on sector size stated above.

The key implication for the purposes of this paper is that, in order to understand the origins of fluctuations, the important shocks are precisely those that affect proportionally more than \sqrt{n} firms.

4 Determining the Number of Factors

The target of estimation in this section will be defined by a complexity parameter τ such that r_1 is the number of factors that affect proportionally more than n^τ cross sections. For the reasons outlined in the previous sections, the number of factors r_1 a practitioner is usually interested in will be such that $\alpha_k > .5$ for $k = 1, \dots, r_1$. This corresponds to complexity parameter $\tau = .5$.

In many applications, the number of factors r_1 is of interest in itself, as illustrated in the last section. For example, we may be interested in the number of fundamental shocks in the economy that contribute to the surprisingly large standard deviation (more than 8 percentage points) of the Federal Reserve Boards Index of Industrial Production (Foerster et al. (2011)). In finance, this number can be interpreted as the number of sources of nondiversifiable risk. In other cases the number of factors must be known to implement various estimation and forecasting procedures. For example, in factor-augmented VAR models, impulse responses based on an incorrect number of factors may be misleading and result in bad policy suggestions (Bernanke et al. (2005), Giannone et al. (2006)). Onatski (2015) discusses the consequences of a misspecified number of factors for the squared error of the estimated common component. I show the implications of the number

of factors on the R^2 of the common component in explaining movements in various series in our empirical section and discuss the effects of the number of factors on squared forecast error in a companion paper (Freyaldenhoven (2017)).

Estimating the number of factors in factor models has been a subject of interest for some time now (e.g. Bai and Ng (2002), Onatski (2010), Ahn and Horenstein (2013)). To the best of my knowledge, all existing estimators are derived from the distribution of eigenvalues of the matrix $\frac{X'X}{T}$ (or equivalently the singular values of X). For example, the information criteria introduced in Bai and Ng (2002) effectively count the number of eigenvalues above a certain threshold, Ahn and Horenstein (2013) consider the ratio of subsequent eigenvalues, while Onatski (2010) uses the difference between subsequent eigenvalues to determine the number of factors.

While the first two methods explicitly require strong factors, “weak” factors are allowed for in Onatski (2010). In the framework of Onatski (2010) some of the “large” eigenvalues do not necessarily diverge to infinity. Essentially, Onatski’s proposed estimator counts the number of eigenvalues which are too large to come from the idiosyncratic errors. While the work of Onatski provide an insightful and novel framework allowing for weak factors, the required assumptions on the error term are quite restrictive. Further, estimating the number of factors from the empirical distribution of eigenvalues still rests on a separability between the two groups of eigenvalues.

In conclusion, all existing methods to estimate the number of factors can be interpreted as formalizing the heuristic approach based upon the visual inspection of the scree plot, which dates back to Cattell (1966).

However, consider the front edge of Figure 2, which depicts the theoretical divergence of an eigenvalue as a function of the corresponding factor strength α . One might consider a thresholding estimator that counts the number of eigenvalues above a threshold $K_n \asymp \sqrt{n}$. Note that the relevant curve is rather flat around this cutoff. This suggests that any thresholding will be very sensitive to the choice of the threshold in finite samples (see, e.g., the discussion in Alessi et al. (2010)). More generally, any estimator trying to distinguish factors on either side of the cutoff based solely on the eigenvalues will share this problem.

I first present an intuitive discussion before I formally state my results. The novel insight here is that in scenarios with local factors, the eigenvectors of the matrix $\frac{X'X}{T}$ carry valuable information which is discarded when solely considering the eigenvalue distribution. I therefore propose to incorporate the eigenvectors into the inference on the number of factors. To this end I introduce the following quantity:

$$\hat{T}_{zk}^u \equiv \psi_k\left(\frac{X'X}{T}\right)\hat{S}_{zk}^u \equiv \psi_k\left(\frac{X'X}{T}\right)\left(\frac{1}{z} \sum_i^z \frac{\hat{\lambda}_{ik}^2}{\sqrt{\frac{1}{n} \sum_{i=1}^n \hat{\lambda}_{ik}^2}}\right)^u, \quad (41)$$

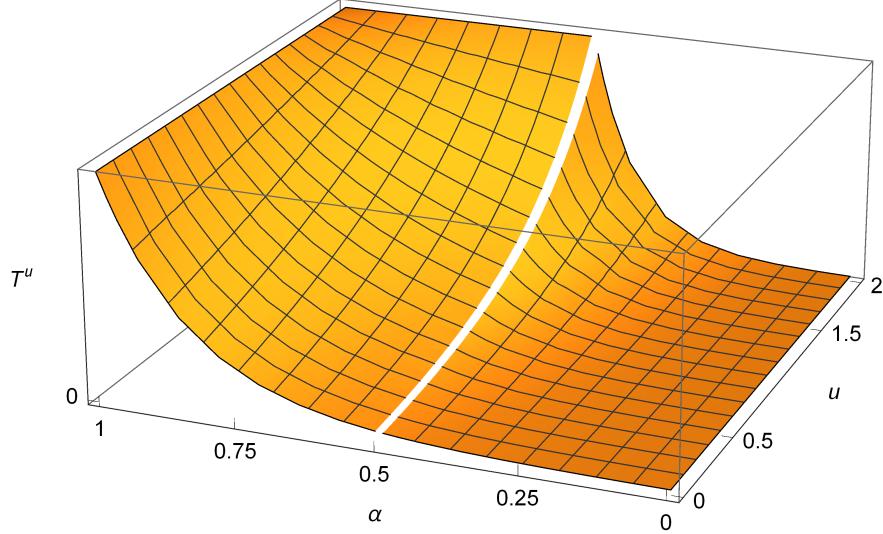


Figure 2: Theoretical divergence rate of \hat{T}_{zk}^u as a function of both factor strength (α) and tuning parameter u for $z = \sqrt{n}$. Note the steep region at $\alpha = \sqrt{n}$ in the back of the picture. $n = 500$.

where, with slight abuse of notation, $\hat{\lambda}_{ik}^2$ are the squared entries of the k th eigenvector sorted in decreasing order, such that I take a partial sum over the z largest entries in the second part. One can think of \hat{T}_{zk}^u as combining the k th eigenvalue of the matrix $\frac{X'X}{T}$ (the first component) with a measure of how concentrated the corresponding eigenvector is on a subset of the observables (the second component). A factor that is highly influential on a subset of covariates but unrelated to the majority of outcomes will be difficult to detect using solely the eigenvalue of the $\frac{X'X}{T}$. However, the second part of (41) will scale this eigenvalue up to enable a practitioner to detect its presence. The power u plays the role of a tuning parameter that governs the relative weight on the eigenvalue versus the eigenvector. With $u = 0$ the second part vanishes and \hat{T}_{zk}^u reduces to just the eigenvalue. On the other hand, with $u = 2$, it can be shown that $\hat{T}_{zk}^2 = \frac{1}{z} \sum_i^z \hat{\lambda}_{ik}^2$ (up to some rescaling). Thus \hat{T}_{zk}^2 only depends on the largest z entries in the k th eigenvector.

This intuition is depicted in Figure 2, which illustrates the behavior of \hat{T}_{zk}^u as a function of factor strength and the tuning parameter u . Consider a DGP with a single factor F . Moving from right to left in Figure 2, the number of covariates affected by the factor increases. On the right edge only a fixed number of covariates is affected, while the left edge of the figure corresponds to a scenario in which the factor is relevant for all covariates. In particular, the number of affected covariates is governed by the parameter $\alpha \in [0, 1]$, with n^α outcomes influenced by the factor. Moving from front to back the value of the tuning parameter u varies from 0 to 2. Note that the position of the steep increase can be chosen by a practitioner through the second tuning parameter z , in this illustration set to $z = \sqrt{n}$. The front edge of the plane, with $u = 0$, simply corresponds to the first eigenvalue of $\frac{X'X}{T}$. With only a fixed number of covariates affected, the first eigenvalue

of the this matrix remains bounded (front right corner). As the factor affects more covariates, the eigenvalue diverges proportionally to the number of covariates that are affected (See Theorem 1). In contrast, on the back edge of Figure 2 the behavior of \hat{T}_{zk}^2 is depicted¹⁰. If a factor affects every outcome or only a fixed number, the behavior of \hat{T}_{zk}^u is invariant to the choice of u (in terms of its rate of divergence). However, exploiting the information contained in the eigenvector, I obtain different divergence rates for factors of intermediate strength.

Suppose we are interested in the number of factors with $\alpha_k > .5$. The flat slope around $\alpha = .5$ on the left edge of Figure 2 suggests that existing estimators, using the eigenvalues of $\frac{X'X}{T}$, will have low discriminatory power in distinguishing factors around this threshold. In contrast, exploiting the information in the eigenvectors (by setting $u > 0$) induces a steep region in the statistic around this threshold.

To derive this result formally, I start by defining the following class of (infeasible) quantities T_{zk}^u . For $u \in [0, 2]$:

$$T_{zk}^u = \psi_k \left(\frac{\Lambda F' F \Lambda'}{T} \right) S_{zk}^u = \psi_k \left(\frac{\Lambda F' F \Lambda'}{T} \right) \left(\frac{1}{z} \sum_i^z \frac{\lambda_{ik}^2}{\sqrt{\frac{1}{n} \sum_{i=1}^n \lambda_{ik}^2}} \right)^u, \quad (42)$$

where λ_{ik}^2 are sorted in decreasing order. Note that for $u = 0$, $T_{zk}^0 = \psi_k(\frac{\Lambda F' F \Lambda'}{T})$.

The behavior of T_{zk}^u is summarized in the following lemma:

Lemma 2. *Under Assumptions 1-2, choose a threshold $z = n^\tau g(n)$, $\tau \in [0, 1]$, such that (i) $g(n) \rightarrow \infty$ and (ii) $g(n)/n^\epsilon \rightarrow 0$ for any $\epsilon > 0$ as $n \rightarrow \infty$. Then, for any given factor $k \leq r$, with $u \in [0, 2]$:*

(a) *If $\alpha_k > \tau$:*

$$T_{zk}^u \asymp n^{(1-\frac{1}{2}u)\alpha_k + \frac{1}{2}u}$$

(b) *If $\alpha_k \leq \tau$:*

$$T_{zk}^u \asymp n^{(1+\frac{1}{2}u)\alpha_k + (\frac{1}{2}-\tau)u} g(n)^{-u}$$

Further, for $k = r+1, \dots, r_{max}$: $T_{zk}^u = 0$

¹⁰Setting $u > 2$ is possible and would result in a quantity that is even more peaked around the threshold parameter τ . The equivalent of Figure 2 extending up to $u = 3$ is depicted in Appendix B. With $u > 2$, T_{zk}^u is no longer monotonically increasing in α_k , the measure of factor strength. I will therefore restrict my analysis to $u \in [0, 2]$ in the remainder of this paper.

Proof. Using Assumption 1 I can rewrite T_{zk}^u as follows:

$$T_{zk}^u = \psi_k \left(\frac{\Lambda F' F \Lambda'}{T} \right) S_{zk}^u = \psi_k \left(\frac{\Lambda F' F \Lambda'}{T} \right) \left(\frac{1}{z} \sum_i^z \frac{\lambda_{ik}^2}{\sqrt{\frac{1}{n} \sum_{i=1}^n \lambda_{ik}^2}} \right)^u \quad (43)$$

$$= \psi_k(\Lambda' \Lambda) \left(\sum_{i=1}^n \lambda_{ik}^2 \right)^{-\frac{1}{2}u} \left(\frac{n^{\frac{1}{2}}}{z} \sum_i^z \lambda_{ik}^2 \right)^u \quad (44)$$

$$= \psi_k(\Lambda' \Lambda)^{1-\frac{1}{2}u} n^{\frac{1}{2}u} \left(\frac{1}{z} \sum_i^z \lambda_{ik}^2 \right)^u. \quad (45)$$

First consider scenario (a). With $\alpha_k > \tau$, the last part of (45) is simply an average of the square of the z largest loadings and therefore

$$T_{zk}^u = O_p(n^{(1-\frac{1}{2}u)\alpha_k}) n^{\frac{1}{2}u} O_p(1). \quad (46)$$

Next, for part (b), let $\alpha_k \leq \tau$: There are only $|\mathcal{A}_k| \asymp n^{\alpha_k}$ non zero loadings in the sum of equation 45. Therefore

$$\frac{1}{z} \sum_i^z \lambda_{ik}^2 = \frac{1}{z} \sum_i^{\mathcal{A}_k} \lambda_{ik}^2 + \frac{1}{z} \sum_i^{\mathcal{A}_k} \lambda_{ik}^2 = O_p\left(\frac{n^{\alpha_k - \tau}}{g(n)}\right) \quad (47)$$

and it follows that

$$T_{zk}^u = O_p(n^{(1-\frac{1}{2}u)\alpha_k}) n^{\frac{1}{2}u} O_p\left(\frac{n^{(\alpha_k - \tau)u}}{g(n)^u}\right). \quad (48)$$

For $k > r$, $\lambda_k = 0$ and the result follows. \square

Since both Λ and F are unobserved, T_{zk}^u is infeasible to compute in practice. I will therefore use the feasible alternative to T_{zk}^u introduced in (41) and repeated below:

$$\hat{T}_{zk}^u \equiv \psi_k\left(\frac{X'X}{T}\right) \hat{S}_{zk}^u \equiv \psi_k\left(\frac{X'X}{T}\right) \left(\frac{1}{z} \sum_i^z \frac{\hat{\lambda}_{ik}^2}{\sqrt{\frac{1}{n} \sum_{i=1}^n \hat{\lambda}_{ik}^2}} \right)^u. \quad (49)$$

Theorem 4. Under Assumptions 1-4, choose a threshold $z = n^\tau g(n)$, $\tau \in [0, 1]$, such that (i) $g(n) \rightarrow \infty$ and (ii) $g(n)/n^\epsilon \rightarrow 0$ for any $\epsilon > 0$ as $n \rightarrow \infty$. Then, for any given Factor $k \leq r_{max}$, with $u \in [0, 2]$:

(a) If $\alpha_k > \tau$:

$$\hat{T}_{zk}^u \asymp n^{\frac{1}{2}u + (1 - \frac{1}{2}u)\alpha_k} \quad (50)$$

(b) If $\max\{\frac{1+\tau}{3}, \frac{\alpha_{max}+4\tau}{6}\} < \alpha_k \leq \tau$:

$$\hat{T}_{zk}^u \asymp n^{(1+\frac{1}{2}u)\alpha_k + (\frac{1}{2}-\tau)u} g(n)^{-u} \quad (51)$$

(c) If $0 < \alpha_k \leq \max\{\frac{1+\tau}{3}, \frac{\alpha_{max}+4\tau}{6}\}$:

$$\hat{T}_{zk}^u \asymp n^{\alpha_k} \quad \text{for } u = 0 \quad (52)$$

$$\hat{T}_{zk}^u = O_p(n^{(1+\frac{1}{2}u)\alpha_k + (\frac{1}{2}-\tau)u} g(n)^{-u}) \quad \text{for } u > 0 \quad (53)$$

(d) If $\alpha_k = 0$ or $k > r$: $\hat{T}_{zk}^u = O_p(1)$

The proof of Theorem 4 can be found in the appendix.

As \hat{T}_{zk}^u is the key quantity in what follows, I also present a corollary that simplifies the notation and covers most cases before I further discuss Theorem 4 and its implications. I argued in section 3 that the important factors are usually those that affect proportionally more than \sqrt{n} of the outcomes, so that $\tau = .5$ will often be the natural choice and I will use this threshold going forward, omitting the corresponding subscript z and writing simply \hat{T}_k^u to obtain the following corollary:

Corollary 1. Let $z = \sqrt{n}g(n)$, such that (i) $g(n) \rightarrow \infty$ and (ii) $g(n)/n^\epsilon \rightarrow 0$ for any $\epsilon > 0$ as $n \rightarrow \infty$. Then, under Assumptions 1-4, for any given Factor $k \leq r_{max}$ and with $u \in [0, 2]$:

(a) If $\alpha_k > \frac{1}{2}$: $\hat{T}_k^u \asymp n^{\frac{1}{2}u + (1-\frac{1}{2}u)\alpha_k}$

(b) If $\alpha_k \leq \frac{1}{2}$:

- $\hat{T}_k^u \asymp n^{\alpha_k}$ for $u = 0$
- $\hat{T}_k^u = O_p(n^{(1+\frac{1}{2}u)\alpha_k} g(n)^{-u})$ for $u > 0$

(c) If $\alpha_k = 0$ or $k > r$: $\hat{T}_k^u = O_p(g(n)^{-u})$

The theoretical rates of Corollary 1 are illustrated graphically in Figure 3. To gain intuition, suppose $\alpha_k = 1$ for $k = 1, 2, \dots, r$, which corresponds to the standard setup in the literature. Then $\hat{T}_{zk}^u \asymp n$ for $k = 1, 2, \dots, r$, regardless of the choice of u (See the left edge of Figure 3). For $k > r$, $\hat{T}_{zk}^0 = O_p(1)$ and $\hat{T}_{zk}^u = O_p(g(n)^{-u})$ if $u > 0$. This means that under the standard setup with only strong factors, the behavior of \hat{T}_{zk}^u is invariant to the choice of u (in terms of its rate of divergence)

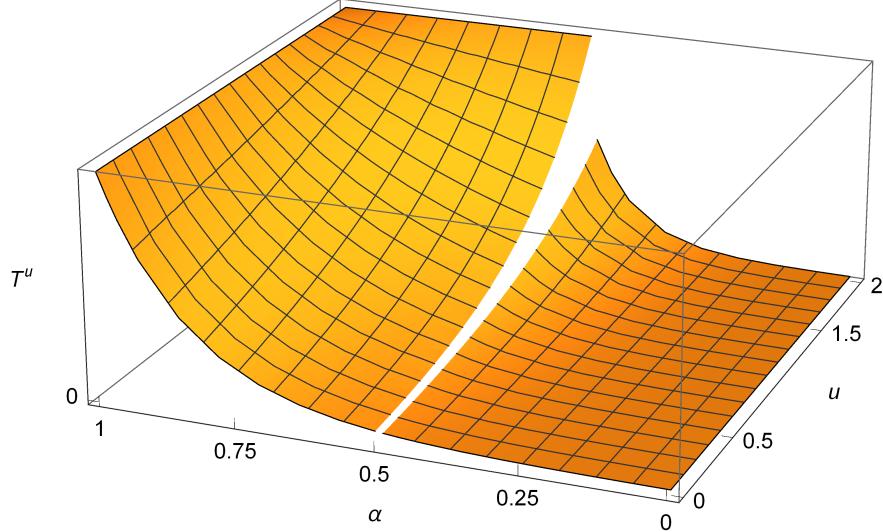


Figure 3: Theoretical rate of \hat{T}_k^u as a function of α_k and u for $z = \sqrt{n} \sqrt{\log\log(n)}$ with $n = 500$

up to the very slowly diverging sequence $g(n)$. In contrast, for all local factors with $\alpha_k \in (0, 1)$ the divergence properties of \hat{T}_{zk}^u depend on the power u . For example, let $\alpha_k \in (.5, 1]$. Then, $\hat{T}_k^0 \asymp n^{\alpha_k}$, $\hat{T}_k^1 \asymp n^{\frac{1}{2} + \frac{1}{2}\alpha_k}$ and $\hat{T}_k^2 \asymp n$. \hat{T}_k^2 has the appealing property that it does not depend on the factor strength for $\alpha_k \in (.5, 1]$. Thus it allows us to distinguish factors above the threshold $\tau = .5$ from idiosyncratic noise at the same rate as strong factors. This property is attractive because in many settings, like those discussed in section 3.2, we are interested in the number of factors r_1 with $\alpha_k > .5$.

An ideal statistic would further be discontinuous with a large jump at a user chosen threshold, thereby making it straightforward to identify the number of factors above this threshold. While \hat{T}_k^2 is discontinuous at τ , the discontinuity is small. However, intuitively the steeper slope leading up to the threshold further suggests better performance for differentiating factors above and below this threshold.

The introduction of the slowly diverging sequence $g(n)$ in the definition of z in Lemma 2 and Theorem 4 will be useful in the construction of estimators in the next section. Specifically, this additional term is responsible for the gap at $\alpha = .5$ and allows us to perfectly separate any factor F_k with factor strength $\alpha_k = .5 + \epsilon$ from a factor F_l with $\alpha_l = .5$ for any $\epsilon > 0$. This relies on the fact that $g(n) = o(n^\epsilon)$ for any $\epsilon > 0$. For most empirically relevant sample sizes, this will only be a good approximation when $g(n)$ diverges extremely slowly. Without the $g(n)$ term, the results will still hold generically, except at the singular point $\alpha_k = .5$, where threshold and divergence rate coincide.

4.1 Proposed Estimators

In this subsection I derive several ways to consistently estimate the number of factors r_1 . I do this by constructing estimators analogous to those that have been proposed in the literature, but using \hat{T}_k^u in place of the empirical distribution of eigenvalues \hat{T}_k^0 . In particular, I focus on the case $u = 2$ and consider the following estimators:

1. An information criteria-like threshold (compare to Bai and Ng (2002), Kapetanios (2004))
2. The difference between two subsequent values (compare to Onatski (2010), Kapetanios (2010))
3. The ratio of two subsequent values (compare to Ahn and Horenstein (2013))

4.1.1 Thresholding Estimators

I start by considering the estimators introduced in Bai and Ng (2002). I will denote by PC the number k that minimizes the criterion function

$$BN(k) = V(k) + k\hat{\sigma}^2 \left(\frac{n+T}{nT} \right) \ln \left(\frac{nT}{n+T} \right), \quad (54)$$

where

$$V(k) = \min_{\Lambda, F_k} (NT)^{-1} \sum_{i=1}^n \sum_{t=1}^T (X_{it} - \lambda_i^k F_t^k)^2 = \frac{1}{nT} \sum_{j=k+1}^n \psi_j(X'X), \quad (55)$$

and $\hat{\sigma}^2$ is an estimator of the unconditional variance of the idiosyncratic error. The second equality follows from the fact that $V(k)$ is the best approximation of X of rank k . We can alternatively represent $\hat{\sigma}^2$ as $V(r_{max}) = \frac{1}{n} \sum_{j=r_{max}}^n \psi_j \left(\frac{X'X}{T} \right)$.¹¹ Therefore, $BN(k)$ is a function of only the empirical distribution of the eigenvalues, and will be equivalent to a thresholding procedure for the aforementioned. Unifying notation in terms of the eigenvalues and using $c = n/T$, this can be seen by rewriting their estimator as:

$$PC = \arg \min_k V(k) + k\hat{\sigma}^2 \left(\frac{n+T}{nT} \right) \ln \left(\frac{nT}{n+T} \right) \quad (56)$$

$$= \arg \min_k \frac{1}{n} \sum_{l=k+1}^n \psi_l \left(\frac{X'X}{T} \right) + k\hat{\sigma}^2 \left(\frac{c+1}{n} \right) \ln \left(\frac{n}{c+1} \right) \quad (57)$$

$$= \max k \quad \text{s.t. } \psi_k \left(\frac{X'X}{T} \right) > \hat{\sigma}^2(c+1) \log \left(\frac{n}{c+1} \right).^{12} \quad (58)$$

¹¹Bai and Ng (2002) consider a total of 6 estimators that differ slightly in their penalty term that is added to $V(k)$ and include a version in logarithms. However, their performances are similar to the ones considered here, and the corresponding results are therefore omitted.

Instead of deriving my estimator solely from the empirical distribution of the eigenvalues, I will consider the following criterion for a fixed constant Q :

$$TC = \max k \quad \text{s.t. } \hat{T}_k^2 > Q \frac{n}{h(n)}, \quad (59)$$

where the function $h(n)$ is such that (i) $h(n) \rightarrow \infty$ and (ii) $h(n)/g(n)^2 \rightarrow 0$ as $n \rightarrow \infty$, and $g(n)$ fulfills the conditions stated in Theorem 4. For example, $h(n) = g(n)$ is a valid choice.

Theorem 5. *Under Assumptions 1-4 TC is a consistent estimator for the number of factors r_1 such that $\alpha_k > .5$ for $k = 1, \dots, r_1$ and $\alpha_k \leq .5$ for $k \geq r_1 + 1$.*

Proof. I first show that $\lim_{n \rightarrow \infty} P(\hat{T}_k^2 > c \frac{n}{h(n)}) = 1$ for $k = 1, \dots, r_1$. In this case, $\alpha_k > .5$. By Theorem 4, $\hat{T}_k^2 \asymp n$. Thus, $n = O_p(\hat{T}_k^2)$. Combining this with $\frac{1}{h(n)} = o_p(1)$ I obtain $\frac{n}{h(n)} = o_p(\hat{T}_k^2)$ and thus

$$\lim_{n \rightarrow \infty} P \left(\left| \frac{\frac{n}{h(n)}}{\hat{T}_k^2} \right| \geq \varepsilon \right) = 0 \quad (60)$$

for any $\varepsilon > 0$. Letting $\varepsilon = 1$ and rearranging I then obtain

$$\lim_{n \rightarrow \infty} P \left(\hat{T}_k^2 \leq \frac{n}{h(n)} \right) = 0. \quad (61)$$

Next, consider the case $\alpha_k \leq .5$. Then $\hat{T}_k^2 = O_p(\frac{n^{2\alpha_k}}{g(n)^2})$ by Theorem 4. But $O_p(n^{2\alpha_k}/g(n)^2) = O_p(n/g(n)^2) = o_p(n/h(n))$ by the definition of $h(n)$ and thus, for any $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} P \left(\left| \frac{\hat{T}_k^2}{\frac{n}{h(n)}} \right| \geq \varepsilon \right) = 0. \quad (62)$$

Letting $\varepsilon = 1$, this completes the proof as I obtain

$$\lim_{n \rightarrow \infty} P \left(\hat{T}_k^2 \geq \frac{n}{h(n)} \right) = 0. \quad (63)$$

□

In practice I propose to incorporate an estimator of the variance into the model, letting $h(n) =$

¹²Similarly Kapetanios (2004) suggests simply using a cutoff value $b = (1 + \sqrt{n/T})^2 + 1$, and estimating the number of factors as the number of empirical eigenvalues above this threshold.

$\frac{g(n)}{Q_2\hat{\sigma}^2}$ such that TC becomes

$$TC = \max k \quad \text{s.t. } \hat{T}_k^2 > Q_1\hat{\sigma}^2 \frac{n}{g(n)} \quad (64)$$

for some fixed constant Q_1 . This is justified because, by Theorem 1

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{j=r_{max}}^n \psi_j \left(\frac{X'X}{T} \right) \leq \frac{1}{n}(n - r_{max} + 1)\psi_{r_{max}} \left(\frac{X'X}{T} \right) \leq \psi_{r_{max}} \left(\frac{X'X}{T} \right) \leq C \quad (65)$$

and, similarly

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{j=r_{max}}^n \psi_j \left(\frac{X'X}{T} \right) \geq \frac{1}{n} \sum_{j=r_{max}}^{[dn]} \psi_j \left(\frac{X'X}{T} \right) \geq c_1 \psi_{[dn]} \left(\frac{X'X}{T} \right) \geq c, \quad (66)$$

where the last inequality uses Weyl's inequalities in a similar way to Theorem 1 to establish that $\psi_{[dn]} \left(\frac{X'X}{T} \right)$ is bounded away from zero, thereby guaranteeing that $\hat{\sigma}^2$ is bounded both above and below. Therefore $h(n) = Qg(n)$, where Q is a finite, positive constant and $h(n)$ fulfills the conditions stated below the definition in (59).

While TC is therefore a consistent estimator for r_1 , I next derive the properties of the existing estimator PC in a setting with local factors. The implicit target of estimation using the PC criterion of Bai and Ng (2002) will be different than the cutoff argued for in this paper. Intuitively, (58) suggests that the PC criterion will estimate the number of factors affecting more than $\log(n)$ outcomes and I indeed obtain the following corollary for the PC estimator:

Corollary 2. *Under Assumptions 1-4 PC is a consistent estimator for the number of factors r^* such that $\alpha_k > 0$ for $k = 1, \dots, r^*$ and $\alpha_k = 0$ for $k \geq r^* + 1$.*

It is therefore clear that PC will not be a consistent estimator for r_1 , unless there exists no Factor k with $0 < \alpha_k \leq 0.5$, in which case r^* and r_1 coincide. However, we can also consider an analogous estimator to PC designed to estimate the number of factors with $\alpha_k > .5$:

$$PC_{\sqrt{n}} = \max k \quad \text{s.t. } \psi_k \left(\frac{X'X}{T} \right) > \hat{\sigma}^2(c + 1) \sqrt{\frac{n}{c + 1}} g(n). \quad (67)$$

It is then easy to show that

Corollary 3. *Under Assumptions 1-4 $PC_{\sqrt{n}}$ is a consistent estimator for the number of factors r_1 such that $\alpha_k > .5$ for $k = 1, \dots, r_1$ and $\alpha_k \leq .5$ for $k \geq r_1 + 1$.*

Given the equivalence established in (58), the proofs of Corollaries 2 and 3 follow the same arguments as the proof of Theorem 5 and are therefore omitted. While this section establishes that

both $PC_{\sqrt{n}}$ and TC are consistent estimators for r_1 , recall the discussion surrounding Figure 3. Based on the steeper slope of \hat{T}_k^2 around the chosen threshold (\sqrt{n}), TC is expected to perform better in finite samples.

4.1.2 Difference Estimators

Instead of choosing a cutoff value, Onatski (2010) establishes that the eigenvalues of the idiosyncratic errors cluster together, while the r eigenvalues corresponding to factors will remain separated. Based on this, one can construct an estimator based on the difference between two subsequent eigenvalues. Starting from r_{max} and successively looking at the difference between two subsequent eigenvalues in decreasing order, the estimator yields $\hat{r} = ED$, the first number at which this difference becomes larger than some constant threshold Q^{13} :

$$ED = \max\{k \leq r_{max} : \psi_k\left(\frac{X'X}{T}\right) - \psi_{k+1}\left(\frac{X'X}{T}\right) \geq Q\} = \max\{k \leq r_{max} : \hat{T}_k^0 - \hat{T}_{k+1}^0 \geq Q\}. \quad (68)$$

Of course, this method also has an analogue when using \hat{T}_k^u for some $u > 0$. Onatski (2010) considers any factors strong enough to be included in the model as soon as their cumulative effects grows with the sample size: r is defined as the number of factors with $\lim_{n \rightarrow \infty} \|\mathcal{A}\| = \infty$. As discussed in previous sections, there are both theoretical and empirical reasons why a practitioner may conclude that some of these factors are too weak to be included in the model. I therefore focus on the case $\tau = .5$ as before and define

$$TD = \max\{k \leq r_{max} : \hat{T}_k^2 - \hat{T}_{k+1}^2 \geq \frac{n}{h(n)}\}, \quad (69)$$

where $h(n)$ is a function such that (i) $h(n) \rightarrow \infty$ and (ii) $h(n)/g(n)^2 \rightarrow 0$ as $n \rightarrow \infty$, and $g(n)$ fulfills the conditions stated in Theorem 4.

Theorem 6. *Under Assumptions 1-4 TD is a consistent estimator for the number of factors r_1 such that $\alpha_k > .5$ for $k = 1, \dots, r_1$ and $\alpha_k \leq .5$ for $k \geq r_1 + 1$.*

Proof. First note that, because $\hat{T}_k^2 = O_p(n^{2\alpha_k}/g(n)^2) = o_p(\frac{n}{h(n)})$ for any k with $\alpha_k \leq 0.5$, $(\hat{T}_k^2 - \hat{T}_{k+1}^2) = o_p(\frac{n}{h(n)})$ for $k > r_1$.

Next consider $k = r_1$. By Theorem 4, if $\alpha_k > .5$, $\lim_{n \rightarrow \infty} P\left(\hat{T}_{r_1}^2 > Q_1 \frac{n}{h(n)}\right) = 1$ and, also by Theorem 4, $\lim_{n \rightarrow \infty} P\left(\hat{T}_{r_1+1}^2 < Q_2 \frac{n}{h(n)}\right) = 1$, for any finite constants $Q_1, Q_2 > 0$. Choosing

¹³Exploiting the shape of the eigenvalue distribution of the idiosyncratic noise at their edge, Onatski (2010) proposes a very appealing way to calibrate the tuning parameter Q that unfortunately is no longer valid in the setup of this paper.

Q_1, Q_2 such that $Q_1 - Q_2 = 1$, gives

$$\lim_{n \rightarrow \infty} P \left((\hat{T}_r^2 - \hat{T}_{r+1}^2) > \frac{n}{h(n)} \right) = 1. \quad (70)$$

□

4.1.3 Ratio Estimators

The most recent estimator that has been introduced to the literature and shown to perform well is based on the ratio of two subsequent eigenvalues following Ahn and Horenstein (2013), defined as

$$ER = \arg \max_{1 \leq k \leq r_{\max}} \frac{\psi_k(\frac{X'X}{T})}{\psi_{k+1}(\frac{X'X}{T})} = \arg \max_{1 \leq k \leq r_{\max}} \frac{\hat{T}_k^0}{\hat{T}_{k+1}^0}. \quad (71)$$

Assumption 5. $\alpha_k > \frac{1}{2}$ for $k = 1, \dots, r_1$ and $\alpha_k = 0$ for $k = r_1, \dots, r_{\max}$.

Because the ratio estimators explicitly relies on a large gap in the eigenvalue distribution, I require an additional assumption of such a gap in Assumption 5 to establish consistency of ratio-based estimators below. Assumption 5 rules out any factors affecting an increasing number of covariates unless the number of affected covariates increases at a rate faster than \sqrt{n} . This assumption is somewhat restrictive, but still less restrictive than the setup of Ahn and Horenstein (2013), who impose $|\mathcal{A}_k| \asymp n$ for $k = 1, \dots, r$. On the other hand the ratio estimator has the significant advantage that it is less dependent on any tuning parameter. It also tends to perform well in finite samples. In line with the ER estimator above I suggest a similar estimator based on the quantity \hat{T}_k^2 :

$$TR = \arg \max_{1 \leq k \leq r_{\max}} \frac{\hat{T}_k^2}{\hat{T}_{k+1}^2}. \quad (72)$$

Theorem 7. Under Assumptions 1-5 TR is a consistent estimator for the number of factors r_1 such that $\alpha_k > .5$ for $k = 1, \dots, r_1$ and $\alpha_k \leq .5$ for $k \geq r_1 + 1$.

Proof. First consider $k = r + 1, \dots, r_{\max}$. Then, by Theorem 4, $\hat{T}_k^2 = O_p(\frac{1}{g(n)^2})$ and thus for every $c_1 > 0$, $\lim_{n \rightarrow \infty} P(\hat{T}_k^2 \geq c_1) = 0$. Further, by Lemma 14 there exists a constant $c_2 > 0$, such that $\lim_{n \rightarrow \infty} P(n\hat{T}_{k^*}^2 \geq c_2) = 1$ for $k^* = r + 1, \dots, r_{\max}$. Then, for any finite $c > 0$, setting

$c_1 = c * c_2$:

$$\begin{aligned} \lim_{n \rightarrow \infty} P\left(\frac{\hat{T}_k^2}{\hat{T}_{k+1}^2} > cn\right) &= \lim_{n \rightarrow \infty} \left[P\left(\frac{\hat{T}_k^2}{\hat{T}_{k+1}^2} > cn \mid \hat{T}_{k+1}^2 < \frac{c_2}{n}\right) P\left(\hat{T}_{k+1}^2 < \frac{c_2}{n}\right) \right. \\ &\quad \left. + P\left(\frac{\hat{T}_k^2}{\hat{T}_{k+1}^2} > cn \mid \hat{T}_{k+1}^2 \geq \frac{c_2}{n}\right) P\left(\hat{T}_{k+1}^2 \geq \frac{c_2}{n}\right) \right] \end{aligned} \quad (73)$$

$$= \lim_{n \rightarrow \infty} P\left(\frac{\hat{T}_k^2}{\hat{T}_{k+1}^2} > cn \mid \hat{T}_{k+1}^2 \geq \frac{c_2}{n}\right) + 0 \quad (74)$$

$$\geq \lim_{n \rightarrow \infty} P(\hat{T}_k^2 > c * c_2) = \lim_{n \rightarrow \infty} P(\hat{T}_k^2 > c_1) = 0. \quad (75)$$

Next, consider $k = r$. By Assumption 5 $\alpha_k > .5$ and there exists a finite $q_1 > 0$ such that $\lim_{n \rightarrow \infty} P(\hat{T}_r^2 > q_1 n) = 1$. Using Assumption 5 again, $\hat{T}_{r+1}^2 = O_p(\frac{1}{g(n)^2})$ and thus for every $q_2 > 0$, $P(\hat{T}_{r+1}^2 \geq q_2) = 0$. Then, for any $q > 0$ and setting $q_2 = q_1/q$:

$$\begin{aligned} \lim_{n \rightarrow \infty} P\left(\frac{\hat{T}_r^2}{\hat{T}_{r+1}^2} > qn\right) &= \lim_{n \rightarrow \infty} \left[P\left(\frac{\hat{T}_r^2}{\hat{T}_{r+1}^2} > qn \mid \hat{T}_{r+1}^2 < q_2\right) P\left(\hat{T}_{r+1}^2 < q_2\right) \right. \\ &\quad \left. + P\left(\frac{\hat{T}_r^2}{\hat{T}_{r+1}^2} > qn \mid \hat{T}_{r+1}^2 \geq q_2\right) P\left(\hat{T}_{r+1}^2 \geq q_2\right) \right] \end{aligned} \quad (76)$$

$$= \lim_{n \rightarrow \infty} \left[P\left(\frac{\hat{T}_r^2}{\hat{T}_{r+1}^2} > qn \mid \hat{T}_{r+1}^2 < q_2\right) + 0 \right] \quad (77)$$

$$\geq \lim_{n \rightarrow \infty} P(\hat{T}_r^2 > q_2 * qn) = \lim_{n \rightarrow \infty} P(\hat{T}_r^2 > q_1 n) = 1. \quad (78)$$

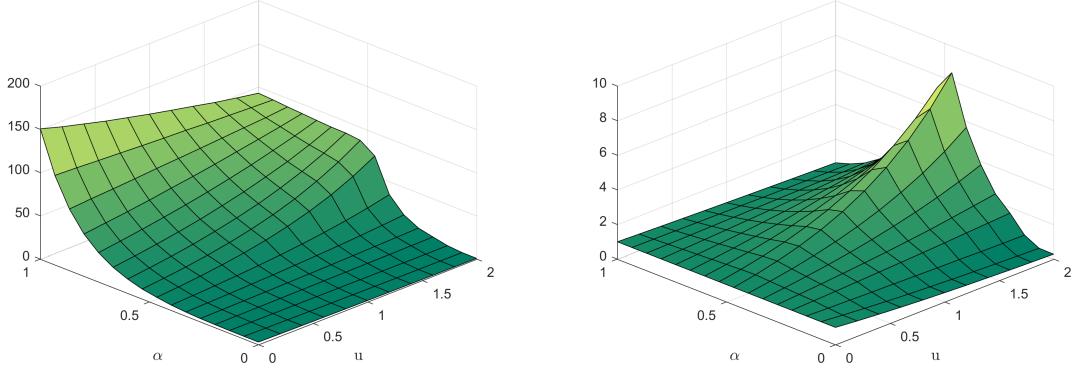
Finally, consider $k = 1, \dots, r-1$. I note that in that case I already established that there exists a finite $q_1 > 0$ such that $\lim_{n \rightarrow \infty} P(\hat{T}_{k+1}^2 > q_1 n) = 1$. It then immediately follows that, for any $c_3 > 0$

$$\lim_{n \rightarrow \infty} P\left(\frac{\hat{T}_k^2}{\hat{T}_{k+1}^2} > c_3 n\right) = 0. \quad (79)$$

□

Corollary 4. *Under Assumptions 1-5 ER is a consistent estimator for the number of factors r_1 .*

The proof largely follows the same arguments as those in the proof of Theorem 7 and is therefore relegated to the appendix. It is perhaps surprising that both estimators require an equal degree of separation. Intuitively, a larger value of u induces a steeper version in \hat{T}_k^u (see Figure 3). As a result, one might expect that the weak and strong factors need to be less well separated to obtain a consistent estimator for higher values of u . To ensure this theoretically a lower bound on \hat{T}_k^2 for



(a) Simulated \hat{T}_{z1}^u as a function of both factor strength α and tuning parameter u
(b) Simulated \hat{S}_{z1}^u as a function of both factor strength α and tuning parameter u

Figure 4: Empirical behavior of proposed quantities. DGP with single factor of varying strength and idiosyncratic noise for $n = 300$, $T = 500$. Figure depicts averages based on 500 simulations.

$k > r_1$, which in turn depends on the partial sum in S^2 , would be needed. I leave this for future research.

On an intuitive level however, the accompanied increase in slope around the targeted factor strength when using \hat{T}_k^u for some $u > 0$ should again improve the performance of this estimator.

5 Simulations

I next present simulation evidence to assess the adequacy of the asymptotic approximations to the finite sample results. In what follows, I fix $g(n) = .7\sqrt{\log\log(n)}$ and note that $g(n)$ clearly fulfills the two criteria stated in Theorem 4: It grows with n , but at a very slow rate and is dominated by n^ε for any $\varepsilon > 0$.¹⁴

I start by depicting the empirical analog to Figure 2 for a simple DGP in Figure 4. The simulated data has a single factor with $T = 500$ and $n = 300$. All loadings are 1 on a subset of covariates with cardinality $|\mathcal{A}| = n^\alpha$ and 0 everywhere else. Error terms u_{it} are i.i.d and standard normal and each variable X_i was centered and divided by its standard deviation before any analysis was performed. In line with Figure 2, α varies from 0 to 1 and u from 0 to 2. Figure 4a depicts the average value of \hat{T}^u across 500 simulations. Note the close resemblance in shape to Figure 2. As u increases, a steep increase in \hat{T}^u emerges around $\tau = .5$. I take Figure 4a as an encouraging sign that the finite sample behavior of \hat{T}^u is well-approximated by the asymptotic theory of Section 3 at least in this simple setting.

¹⁴I chose $g(n) = .7\sqrt{\log\log(n)}$ because it is approximately equal to one for most relevant sample sizes.

I also emphasize that of the two quantities depicted on the horizontal axes, α is unknown to a practitioner, while u is a tuning parameter that can be varied. This suggests u can be varied as an exploratory tool. If F_k is a local factor in the sense of this paper, the divergence rate of \hat{T}_k^u changes as u increases. Although not explicitly in my model, it is clear that the same does not hold if F_k is a weak, but global factor (i.e. if F_k has a small effect on all outcomes). By using only the eigenvalues of $\frac{X'X}{T}$, factors with a strong effect on a subset of outcome and factors with a weak effect on all outcomes will be treated equally. However, depending on the economic model or context, a researcher may be more interested in one or the other. The change in shape associated with an increase in the tuning parameter u is therefore indicative of the underlying structure and a practitioner might be interested in the behavior of \hat{T}_z^u when u increases from 0 to 2. Since $\hat{S}_k^u = \hat{T}_k^u / \hat{T}_k^0$, this amounts to looking at \hat{S}_k^2 (the “peakedness” of the eigenvector) directly. With $\tau = \frac{1}{2}$ for simplicity, this quantity behaves as follows:

$$\hat{S}_k^2 = \hat{T}_k^2 / \hat{T}_k^0 \asymp \begin{cases} n^{1-\alpha_k} & \text{for } \alpha_k \geq \frac{1}{2} \\ O_p(n^{\alpha_k} / g(n)^2) & \text{for } \alpha_k < \frac{1}{2}. \end{cases} \quad (80)$$

For the simple DGP introduced above, \hat{S}_1^u is depicted in Figure 4b. It suggests that empirical data approximates this quantity rather well at least in some simple DGPs: While the eigenvalue is monotonically increasing in factor strength, \hat{S}_k^2 takes its highest value at $\alpha_k = .5$.

I next consider more realistic settings as they might be observed in practice. I consider a panel with

$$\underset{(500 \times 300)}{X} = \underset{(500 \times 6)(6 \times 300)}{F} \underset{(500 \times 6)}{\Lambda'} + \underset{(500 \times 3)(3 \times 300)}{G} \underset{(500 \times 3)}{\Lambda^{w'}} + \sqrt{\theta} \underset{(500 \times 300)}{e}, \quad (81)$$

where $(T, n) = (500, 300)$ falls within the range of dimensions usually considered in the literature¹⁵ and will be varied later on. The variables exhibit a factor structure with 6 independent factors $F_k, k = 1, 2, \dots, 6$, drawn from a standard normal distribution. The 500×6 loading matrix Λ is created by filling random subsets of its columns with $(1 + \eta_{ik})$, where η_{ik} is drawn from a standard normal. These subsets will be of varying size and dictate which variables are affected by the corresponding factor, with the sequence of group sizes given by $\{|\mathcal{A}_k|\}_{k=1}^6 = \{n, n^{.85}, n^{.75}, n^{2/3}, n^{2/3}, n^6\}$ rounded to the nearest integer for the 6 Factors. All other entries in Λ are zero. There are three additional factors G_1, G_2, G_3 also drawn from a standard normal, which I consider too weak to be pervasive. Their loading matrix Λ^w has entries $(1 + \eta_i)$, where η_i is drawn from a standard normal on random subsets of its columns with cardinalities $n^{1/3}, n^{1/4}$ and $\log(n)$, again rounded to the nearest integer. All remaining entries are zero. For the idiosyncratic part I allow for both cross sectional and inter temporal correlation. I follow Onatski (2010),

¹⁵For example Bai and Ng (2002) consider sample sizes in both dimensions between 40 and 8000.

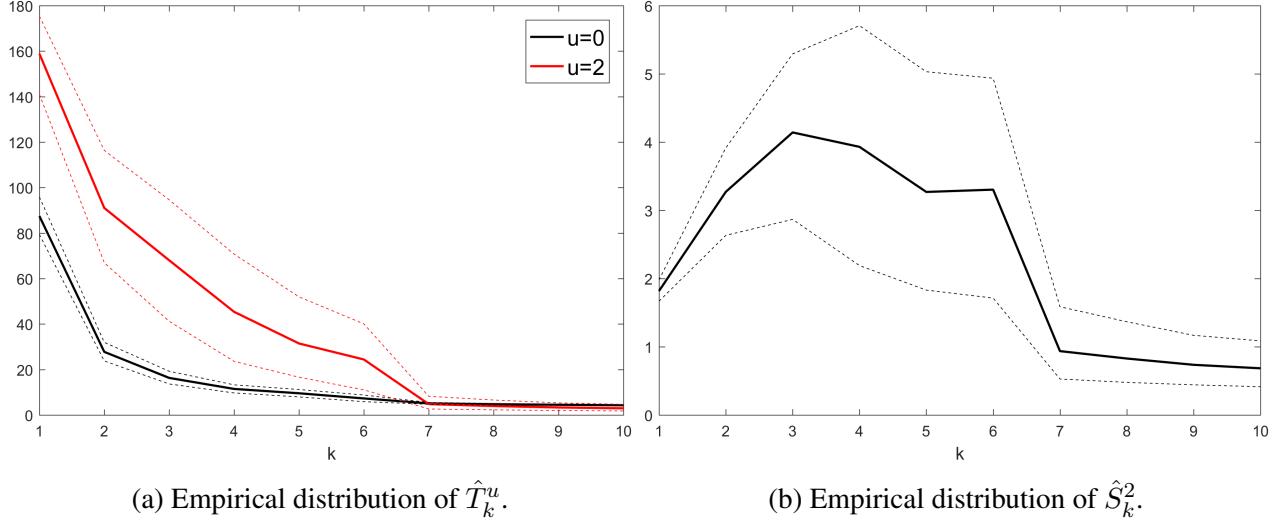


Figure 5: Illustration of key quantities in a simulated dataset. Solid line depicts average across 1000 realizations. dashed lines illustrate 5th and 95th quantile. Data generated by baseline DGP, with $n = 300$, $T = 500$, $(\rho, \beta) = (0.3, 0.1)$, $r = 6$, $\theta = 1.5$. x-axis corresponds to first ten eigenvalues/eigenvectors.

Onatski (2012) and model the errors as follows:

$$e_{ti} = \rho e_{t-1,i} + (1 - \rho^2)^{1/2} v_{it} \quad (82)$$

$$v_{ti} = \beta v_{t,i-1} + (1 - \beta^2)^{1/2} u_{it}, \quad u_{it} \sim i.i.d N(0, 1), \quad (83)$$

with baseline parameter values of $(\rho, \beta) = (0.3, 0.1)$ as in Onatski (2010). Finally I can think of the constant θ as varying the signal-to-noise ratio and set $\theta = 1.5$ in my baseline model. The factor structure and signal-to-noise ratio of the baseline DGP are designed to closely reproduce the screeplot in the macroeconomic application (see Figure 1 in the introduction).

Figure 5 depicts the behavior of both \hat{T}_k^u , $u \in \{0, 2\}$ and \hat{S}_k^2 in the data. It is constructed from 1000 realizations, with the bold line depicting the average and the dashed lines depicting the pointwise 5th and 95th percentile of the respective quantities. Noting that \hat{T}_k^0 and \hat{T}_k^2 in Figure 5a correspond to the left and right edge of Figure 2 respectively, I observe an encouraging resemblance with a larger jump at $\hat{r} = 6$ when $u = 2$. This is due to the behavior of S_k^u , depicted in Figure 5b. The eigenvectors corresponding to more local factors are indeed more concentrated on a subset of its entries.

I next depict the ratios and differences of subsequent values of \hat{T}_k^u in Figure 6. Consider an estimator constructed as the maximum of subsequent ratios of \hat{T}_k^u , which are depicted in Figure 6a. In contrast to an estimator derived solely from the eigenvalues of $\frac{X'X}{T}$ (ER), which suggests the presence of a single factor based on the average depicted here, incorporating the eigenvectors by

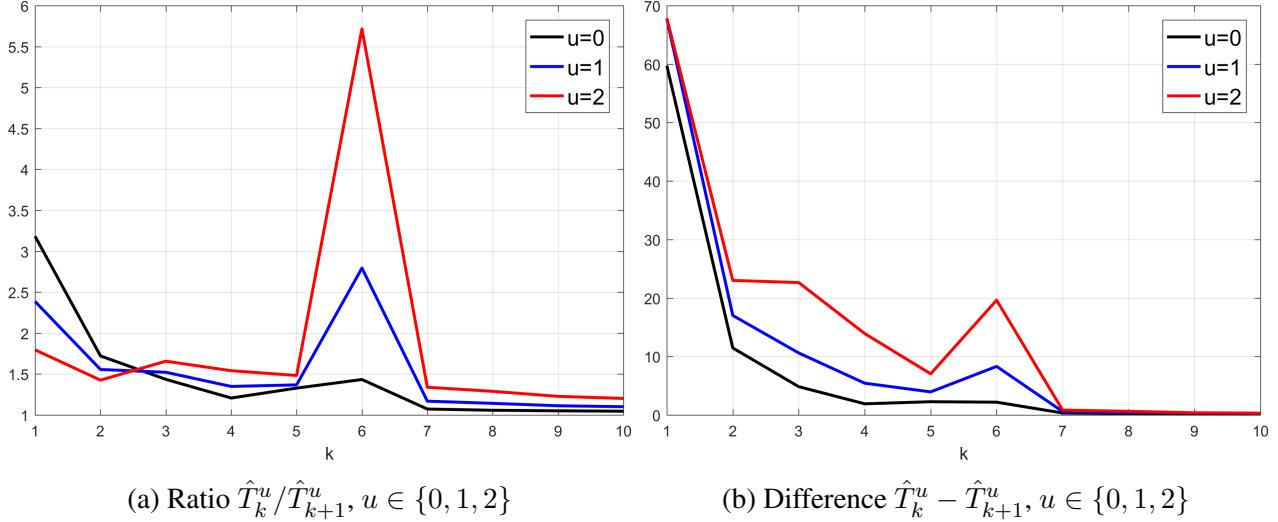


Figure 6: Depicted are averages across 1000 replications in a simulated dataset. Data generated by baseline DGP, with $n = 300$, $T = 500$, $(\rho, \beta) = (0.3, 0.1)$, $\theta = 1.5$, $r = 6$. x-axis corresponds to first ten eigenvalues/eigenvectors

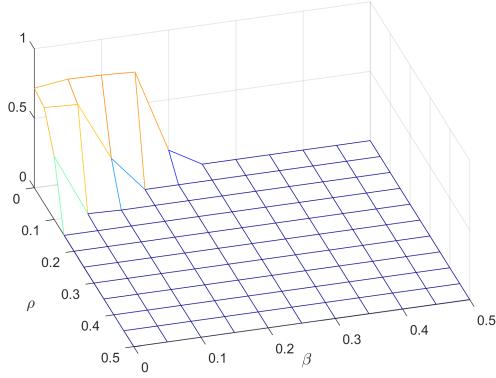
setting $u = 2$ (TR) clearly yields an estimate of $\hat{r} = 6$. For an estimator based on the differences (Figure 6b) it is more difficult to gage what the estimator would select from the picture, but we similarly observe a larger jump at $k = 6$ as the tuning parameter u increases.

In the remainder of this section I will demonstrate how the various estimators from section 4.1 perform for varying amounts of correlation in the error terms, various values of the signal-to-noise ratio, and varying sample sizes. The baseline DGP is the one in (81) and described thereafter.

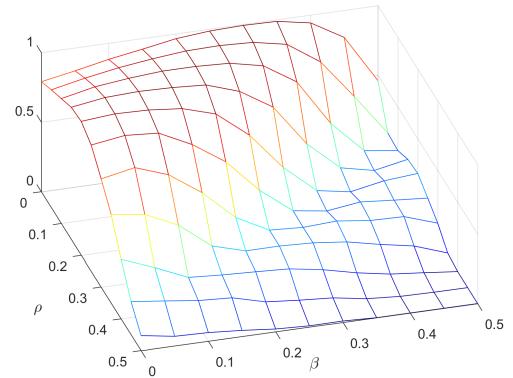
Figures 7 - 8 depict the performance of the estimators for the number of factors as the correlation in the idiosyncratic noise increases. In these figures, I vary the dependence structure of the error term along a two dimensional grid from idiosyncratic noise to include significant cross-sectional dependence and autocorrelation. Where applicable, I let $h(n) = .1\hat{\sigma}^2\sqrt{\log\log(n)}$.

Figure 7 depicts the percentage of simulations in which an estimator correctly estimates the number of factors to be 6. Figure 8 depicts the average number of factors an estimator yields across simulations. The first panel uses the PC criterion of Bai and Ng (2002) to estimate the number of factors. In Panels 7b and 8b I use the thresholding estimator TC . The second row uses the maximum ratio of two subsequent values of $\psi_k(\frac{X'X}{T})$ and T_k^2 respectively, corresponding to estimators ER (Ahn and Horenstein (2013)) and TR . The left and right side of the first two rows are therefore directly comparable: The left panels depict the results of the existing estimators based on the eigenvalues, while the right panels depict the corresponding estimators that incorporate the information in the eigenvector.

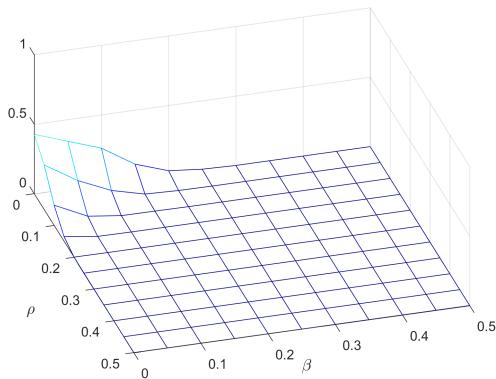
The bottom left panels depict the estimated number \hat{r} using the ED estimator of Onatski (2010), while the bottom right implements an alternative thresholding estimator based on only



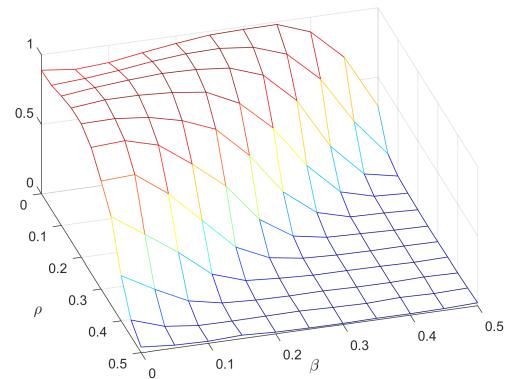
(a) Information criterion PC from Bai and Ng (2002) (PC)



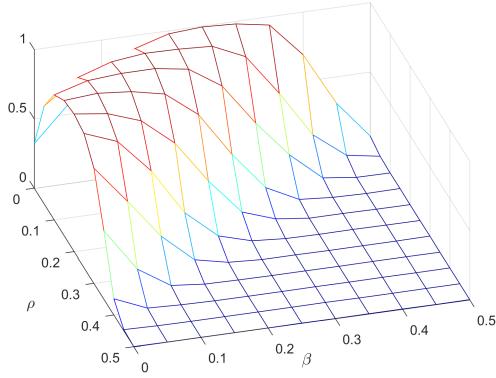
(b) Thresholding based on T^2 (TC)



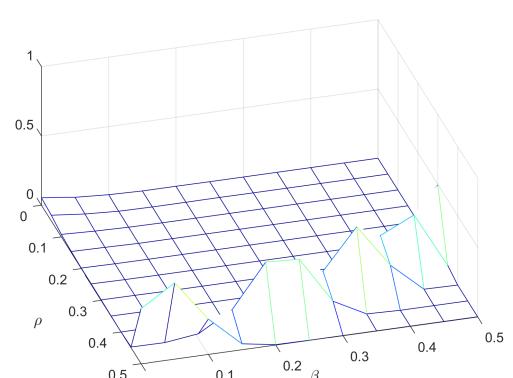
(c) Maximum ratio of two subsequent eigenvalues as in Ahn and Horenstein (2013) (ER)



(d) Maximum ratio of two subsequent values of T^2 (TR)

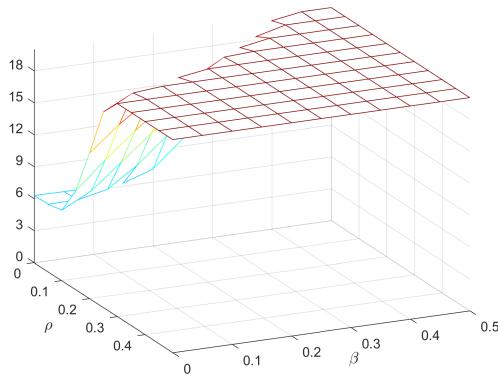


(e) Difference of two subsequent eigenvalues as in Onatski (2010) (ED)

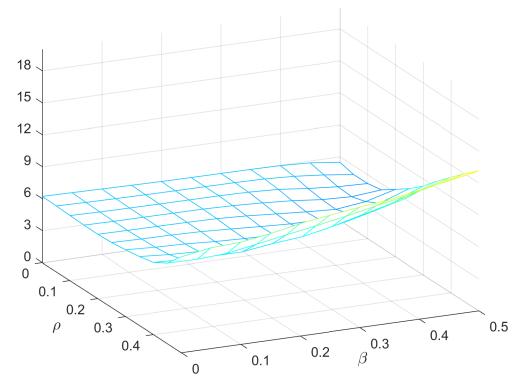


(f) Thresholding based on T^0 (PC $_{\sqrt{n}}$)

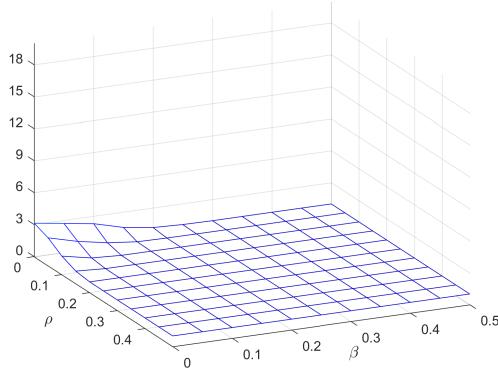
Figure 7: Percentage of simulations correctly estimating the number of “relevant” factors r_1 as both cross sectional and intertemporal correlation is varied. $r_1 = 6$ and $r = 9$, with later factors affecting fewer cross sections ($\{\mathcal{A}_k\}_{k=1}^9 = \{n, n^{.85}, n^{.75}, n^{2/3}, n^{2/3}, n^{.6}, n^{1/3}, n^{1/4}, \log(n)\}$). $n = 300$, $T = 500$, $\theta = 1.5$. Figure based on 500 replications.



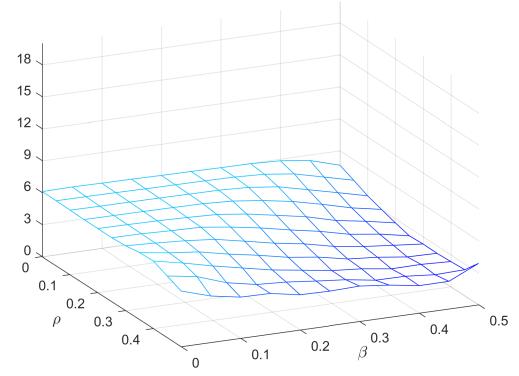
(a) PC criterion from Bai and Ng (2002) (PC)



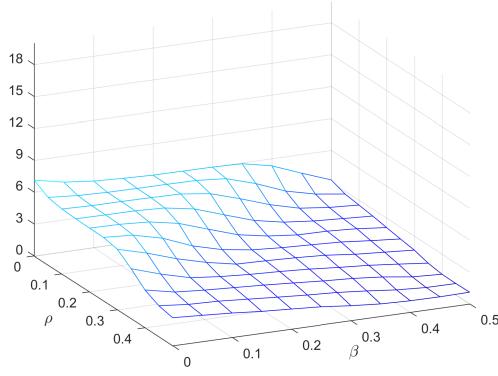
(b) Thresholding based on T^2 (TC)



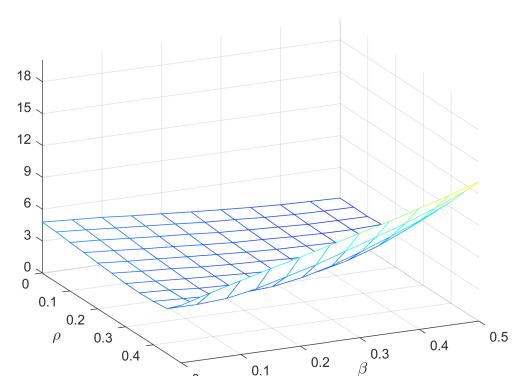
(c) Maximum ratio of two subsequent eigenvalues as in Ahn and Horenstein (2013) (ER)



(d) Maximum ratio of two subsequent values of T^2 (TR)



(e) Difference of two subsequent eigenvalues as in Onatski (2010) (ED)



(f) Thresholding based on T^0 ($PC_{\sqrt{n}}$)

Figure 8: Average estimate for number of factors as both cross sectional and intertemporal correlation is varied. $r_1 = 6$ and $r = 9$, with later factors affecting fewer cross sections ($\{|\mathcal{A}_k|\}_{k=1}^9 = \{n, n^{.85}, n^{.75}, n^{2/3}, n^{2/3}, n^{.6}, n^{1/3}, n^{1/4}, \log(n)\}$). $n = 300$, $T = 500$, $\theta = 1.5$. Figure based on 500 replications.

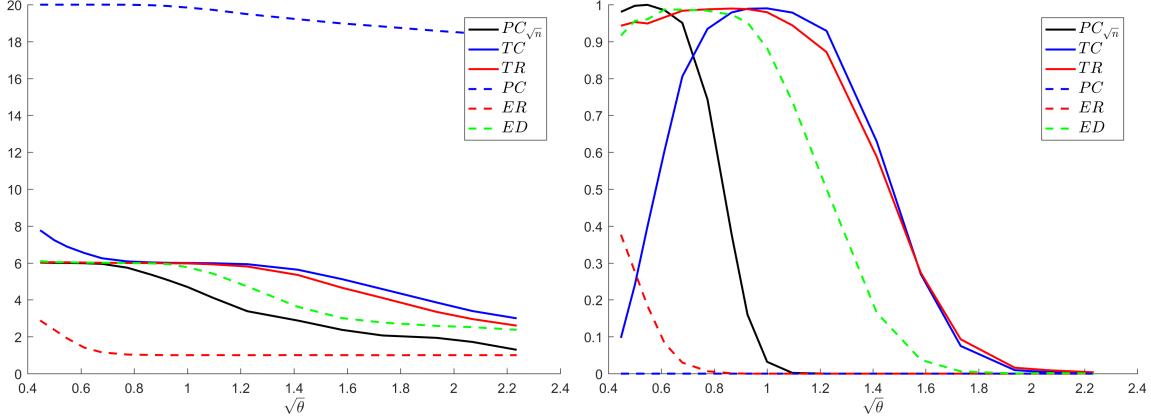
the eigenvalues in $PC_{\sqrt{n}}$.

I conclude from Figure 7 that, even under an exact factor structure (i.e. with uncorrelated errors across both cross sections and time), only the two estimators incorporating the additional information in the eigenvectors (Figures 7b and 7d) come close to correctly estimating $r = 6$ across simulations. In fact, three of the estimators I consider here fail to correctly detect the number of factor in more than 90% of my simulations anywhere the entire parameter space. Of the estimators derived from the eigenvalues, the difference based estimator of Onatski (2010) performs best. As outlined in the previous section, both the PC and ED estimator are designed to also detect very weak factors. This would lead them to overestimate the number of factors here due to the presence of the three additional factors G . However, based on Figure 8, the PC of Bai and Ng (2002) on average yields an estimate \hat{r}^{PC} close to 6, while it is clear from 7 that there is substantial variation left in those estimates. This reenforces my previous argument that any thresholding based on the eigenvalue distribution will be imprecise. Further, as the correlation in the errors increases, the estimate based on the PC estimator deteriorates and quickly approaches $r_{max} = 20$. The ED estimator selects on average between 7 and 8 Factors under a strict factor model. However, as the amount of correlation in the error term increases, the number of estimated factors decreases and, consequently $\hat{r}^{ED} = 6$ for modest values of correlation before it begins to underestimate the number of factors.

In contrast, the estimators based upon T^2 are more robust to the introduction of dependence in the errors. They are on average correct in the simple case of no correlation in the error terms and remain correct for modest levels of cross-sectional and inter-temporal correlation. In particular, when comparing the left and right hand side of the first two rows, where the estimators are directly comparable, we observe a clear benefit from setting $u > 0$. Based on Figures 7 and 8 I conclude that the TR estimator, based on the ratio of subsequent values of \hat{T}_k^2 performs best across the parameter space. Further, Onatski (2010) argues that the parameter pair $(\rho, \beta) = (0.3, 0.1)$ describes the data well in many financial applications. At those parameter values, my Monte Carlo simulations point to a significant gain in performance. Specifically, on average $\hat{r}^{TR} = 5.83$ while the best performance of any of the existing estimators is again ED with $\hat{r}^{ED} = 4.51$.

I will next vary the signal-to-noise ratio in my data by varying θ , while fixing the correlation structure in the errors at $(\rho, \beta) = (0.3, 0.1)$. Figure 9 depicts the sensitivity of the estimators to the amount of noise in the data. I again conclude that the information inherent in the eigenvectors greatly improves the estimation of the number of factors. My proposed estimators remain correct for a considerably larger range of θ compared to their counterparts derived solely from the empirical eigenvalues. Of those, the ED estimator of Onatski (2010) again performs best. Nevertheless, I again conclude that the TR estimator performs best in this setting.

For the final set of simulations I vary the dimension of the dataset by varying both the cross-



(a) Average estimate for number of factors according to TC , TR , $PC_{\sqrt{n}}$, PC , ER and ED (b) % of simulations correctly estimating number of factors r_1 according to TC , TR , $PC_{\sqrt{n}}$, PC , ER and ED

Figure 9: Empirical behavior of estimators as the relative variance of idiosyncratic noise increases. $r_1 = 6$ and $r = 9$, with later factors affecting fewer cross sections ($\{|\mathcal{A}_k|\}_{k=1}^9 = \{n, n^{.85}, n^{.75}, n^{2/3}, n^{2/3}, n^{.6}, n^{1/3}, n^{1/4}, \log(n)\}$). $n = 300$, $T = 500$, $(\rho, \beta) = (0.3, 0.1)$. Figure based on 500 replications.

sectional dimension as the well as the time horizon of the data. Table 2 depicts the results. I again fix $(\rho, \beta) = (0.3, 0.1)$ and $\theta = 1.5$. Each entry in Table 2 consists of two numbers $\hat{r}/\%$, where \hat{r} is the average number of estimated factors, and $\%$ is the percentage correctly classifying $r = 6$. In small samples all estimators perform poorly. While the ratio- and difference-based estimators tend to underestimate the true number of factors, both estimators based on thresholding the eigenvalues tend to severely overestimate the number of factors in small samples. Again comparing ER and TR , as well as PC and TC directly, the previous pattern holds up: Setting $u > 0$ significantly improves the performance of the estimator. Within the class of estimators derived from solely the eigenvalues, the ED estimator of Onatski (2010) tends to perform best. However, it is interesting to observe that, similar to the pattern observed in Figure 9, the ED estimator consistently underestimates the number of factors, even though it is constructed explicitly to allow for and detect weaker factors in the data. Especially as the sample size increases, the TR estimator tends to perform best across the estimators considered.

In conclusion I find that the TR estimator tends to perform best across most of the DGPs considered here, but note that it tends to underestimate the number of factors in small sample sizes, where no estimator does particularly well.

In the appendix I repeat the analysis of this section with an alternative DGP that has a strong factor structure. In particular, I use the same baseline DGP as in this section but set $\lambda_{ik} = 1 + \eta_{ik}$, where η_i is still drawn from a standard normal, for every entry in Λ and exclude the very weak factors G . Thus $\alpha_k = 1$ for $k = 1, \dots, r$ as is usually the case in the literature. I find that, under

<i>n</i>	<i>T</i>	<i>ER</i>	<i>TR</i>	<i>PC</i>	<i>PC</i> _{\sqrt{n}}	<i>TC</i>	<i>ED</i>
100	100	1.03 / 0.00	4.02 / 0.02	20 / 0.00	14.7 / 0.00	14.3 / 0.00	1.53 / 0.00
100	150	1.02 / 0.00	3.52 / 0.03	20 / 0.00	11.8 / 0.00	12.4 / 0.00	1.74 / 0.00
150	100	1.02 / 0.00	3.16 / 0.02	20 / 0.00	13.7 / 0.00	14.6 / 0.00	1.7 / 0.00
150	250	1.01 / 0.00	3.86 / 0.09	20 / 0.00	5.94 / 0.61	8.62 / 0.04	2.2 / 0.02
150	500	1.01 / 0.00	4.98 / 0.42	19.1 / 0.00	4.6 / 0.05	6.25 / 0.62	3.28 / 0.20
300	250	1 / 0.00	4.39 / 0.19	20 / 0.00	4.09 / 0.00	6.5 / 0.45	2.72 / 0.02
300	500	1 / 0.00	5.82 / 0.88	19.5 / 0.00	3.38 / 0.00	5.93 / 0.93	4.62 / 0.49
300	750	1 / 0.00	5.96 / 0.97	16.9 / 0.00	3.28 / 0.00	5.97 / 0.97	5.76 / 0.86
500	250	1 / 0.00	4.58 / 0.16	20 / 0.00	3.16 / 0.00	5.81 / 0.45	3.02 / 0.01
500	500	1 / 0.00	5.92 / 0.94	20 / 0.00	3.01 / 0.00	5.95 / 0.95	5.27 / 0.62
500	750	1 / 0.00	6 / 1.00	17.8 / 0.00	3 / 0.00	5.99 / 0.99	5.95 / 0.95
1000	1000	1 / 0.00	6 / 1.00	16.8 / 0.00	2.98 / 0.00	6 / 1.00	6 / 1.00

Table 2: Each entry depicts a combination $\hat{r}/\%$, where \hat{r} is the average number of estimated factors, and $\%$ is the percentage correctly classifying $r_1 = 6$. In each row, the highest percentage is highlighted. $r_1 = 6$ and $r = 9$, with later factors affecting fewer cross sections ($\{\mathcal{A}_k\}_{k=1}^9 = \{n, n^{.85}, n^{.75}, n^{2/3}, n^{2/3}, n^{.6}, n^{1/3}, n^{1/4}, \log(n)\}$). $n = 300, T = 500, \theta = 1.5$. Table based on 500 replications.

a strong factor structure, estimators incorporating the partial sums in the eigenvector generally perform no worse than existing estimators, although the *ED* estimator of Onatski (2010) tends to perform particularly well in smaller samples. I therefore conclude that raising \hat{T}_k^u to a power $u > 0$ has little implications if all factors are strong, but yields significant performance gains if local factors are present in the data.

In conclusion, my recommendation for estimating the number factors r_1 that affect proportionally more than \sqrt{n} variables is therefore to use the *TR* estimator with its implementation outlined as follows:

1. Obtain preliminary estimates $\hat{F}, \hat{\Lambda}$ using the first r_{max} principal components, where r_{max} is large enough such that $\psi_k(\frac{X'X}{T})$ is bounded for $k > r_{max}$.
2. Letting $z = 0.7\sqrt{\log(\log(n))}\sqrt{n}$, rounded to the nearest integer, compute

$$\hat{T}_{zk}^2 \equiv \psi_k(\frac{X'X}{T})\hat{S}_{zk}^u \equiv \psi_k(\frac{X'X}{T})\left(\frac{1}{z} \sum_i^z \frac{\hat{\lambda}_{ik}^2}{\sqrt{\frac{1}{n} \sum_{i=1}^n \hat{\lambda}_{ik}^2}}\right)^2. \quad (84)$$

Note that $g(n) = .7\sqrt{\log\log(n)} \approx 1$ for most relevant sample sizes, and this recommendation is therefore generic.

3. Set

$$\hat{r} = TR = \arg \max_{1 \leq k \leq r_{\max}} \frac{\hat{T}_k^2}{\hat{T}_{k+1}^2}. \quad (85)$$

6 A factor model with local factors of the US economy

Two classic applications where factor models have proven particularly useful are macroeconomic monitoring and forecasting (see Stock and Watson (2016) for a good review). This section describes the factor model estimated from a large panel of US macroeconomic indicators under the weaker assumptions maintained in this paper and illustrates the implications of the presence of local factors.

I employ one of the standard datasets in the factor model literature in macroeconomics (see, e.g. Stock and Watson (2005) and De Mol et al. (2008)). The data contains quarterly observations of 207 macroeconomic variables, primarily for the US economy. In particular, I use the vintage of the dataset used in the handbook chapter of Stock and Watson (2016). It includes real activity variables, prices, productivity and earnings, interest rates and spreads, money and credit, asset and wealth variables, oil market variables and indicators representing international activity. The data ranges from 1959Q1-2014Q4. All variables have been transformed to achieve approximate stationarity and a small number of outliers were removed. I follow the same transformations as Stock and Watson (2016) and also follow their practice in removing low-frequency trends in the data using a biweight low-pass filter, with a bandwidth of 100 quarters, as in Stock and Watson (2012)¹⁶.

I note that the dataset consists of series at multiple levels of aggregation and therefore some of the series are by construction closely related to linear combinations of other series. (Although, since series are transformed to achieve stationarity, these aggregation identities no longer hold exactly.) For example, the dataset not only includes seven sectoral measures of industrial production, but also three aggregate measures of industrial production, constructed from the former seven. I only use the disaggregated time series in my estimation of the factor structure and disregard the aggregates. The reasoning for this is outlined in the Appendix. For a more detailed discussion of this issue, also see Boivin and Ng (2006). This elimination leaves 139 variables in the data.

Only 94 of those series are available for the entire sample and I will restrict my analysis to the 94 time series that have availability for the entire sample. This allows for a straightforward implementation of the principal component estimator.¹⁷

¹⁶Data are available at <http://www.princeton.edu/~mwatson/publi.html>. For a full description of the data, as well as a more detailed description of the transformations to the raw data I refer the reader to Stock and Watson (2016).

¹⁷Alternatively one could analyze the full sample of 139 disaggregated variables using the EM algorithm of Stock

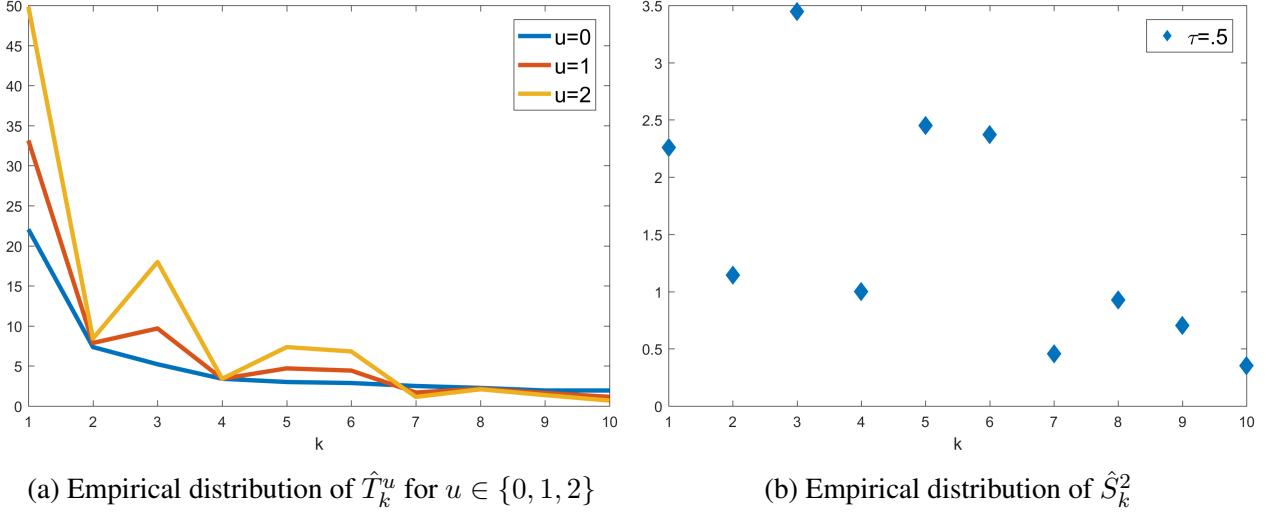


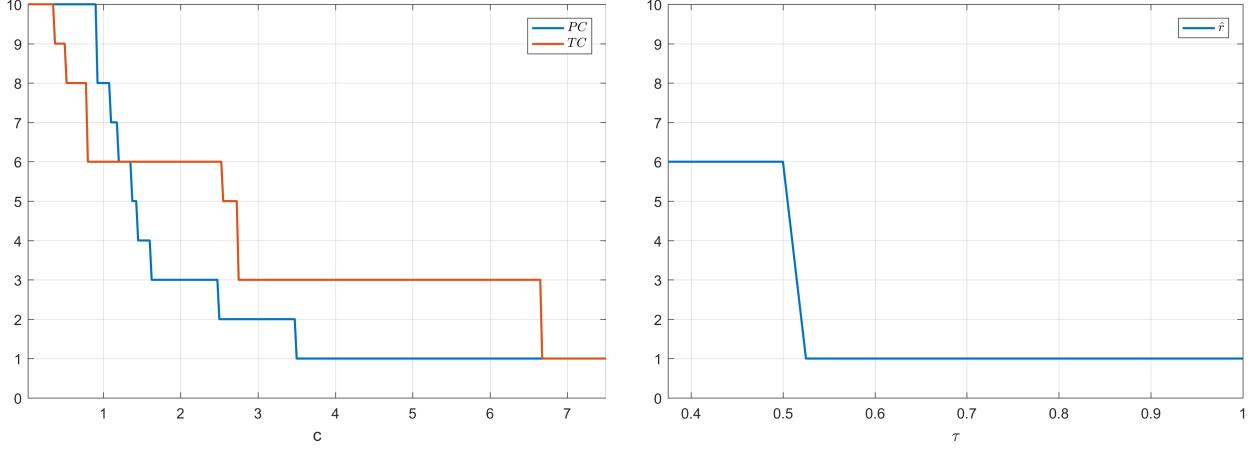
Figure 10: Illustration of key quantities for $z = \sqrt{n}$ in a dataset of macroeconomic indicators in the US

I start by depicting some of the key quantities used throughout the paper in this dataset and provide some intuitive discussion of those. Figure 10 depicts the behavior of both \hat{T}_k^u , $u \in \{0, 1, 2\}$ and \hat{S}_k^2 in the data (Note that setting $u = 0$ reproduces Figure 1 from the introduction). I note that incorporating information from the eigenvector (setting $u > 0$) means that the ordering of \hat{T}_k^u may be reversed. If \hat{T}_k^u corresponds to a factor with a strong local effect—heavily influencing a subset of outcomes—it will be scaled up significantly if $u > 0$. On the other hand, this does not hold for an eigenvalue corresponding to a factor that weakly affects all outcomes. This is illustrated in Figure 10b: Some of the eigenvectors are more concentrated on a subset of its entries than others. In particular, with a tuning parameter $\tau = .5$, I find that the 3rd, 5th and 6th eigenvector are particularly concentrated on its largest loadings. As a consequence we see in Figure 10a that the order of \hat{T}_k^u changes as u increases (e.g. \hat{T}_5^2 is larger than \hat{T}_4^2). A visual inspection of Figure 10a also indicates a drop-off at $k = 3$ and $k = 6$, suggesting the presence of either 3 or 6 factors, depending on the minimum strength of the factors a practitioner would like to keep in her model.

I next summarize the results of the 6 estimators considered throughout this paper in Table 3. While both estimators derived from \hat{T}_k^2 suggest the presence of 6 factors in the data, the three existing estimators from the literature considered here (*ER*, *PC*, *ED*) find evidence for 1, 3 and 8 and Watson (2002b) to handle missing observations.

Estimator	<i>ER</i>	<i>TR</i>	<i>PC</i>	<i>PC</i> _{\sqrt{n}}	<i>TC</i>	<i>ED</i>
Estimated number of factors	1	6	8	3	6	3

Table 3: Estimated number of factors for the 6 estimators considered throughout this paper.



(a) Estimated number of factors by varying the cut-off (cf. Alessi et al. (2010)). Cutoffs from Theorem 5 and Corollary 2 multiplied by constant c .

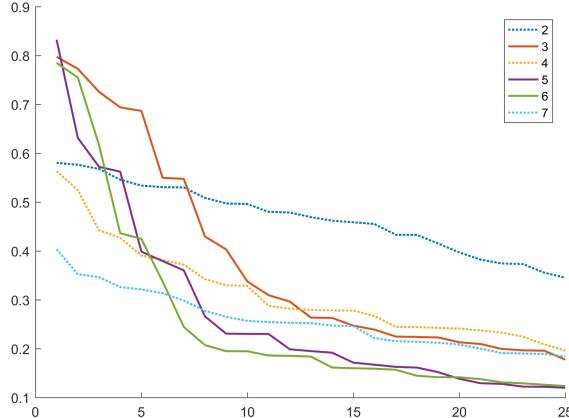
(b) TR , the maximum of $\hat{T}_k^2/\hat{T}_{k+1}^2$, $k \in \{2, 3, \dots, 10\}$ for varying tuning parameter $\tau \in [.375, 1]$.

Figure 11: Illustration of estimators in dataset of US macroeconomic indicators for varying tuning parameters

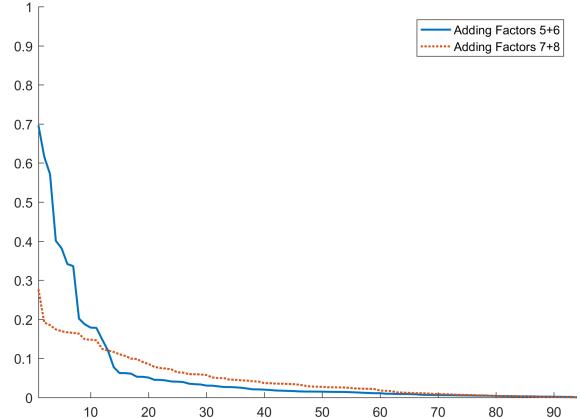
factors respectively.

To address the problem that the estimated number of factors is sensitive to the choice of cutoff under the thresholding estimator, Alessi et al. (2010) suggests to vary this threshold and explore how the estimated number of factors changes. The result is depicted in Figure 11a. It depicts the estimated number of factors based on PC and TC as a function of the tuning parameter c , which multiplies the threshold (Noting that the thresholds in Theorem 5 and Corollary 2 can be arbitrarily scaled by a constant). The discussion in section 4 suggests that incorporating the eigenvector should make the estimator less dependent on the choice of the tuning parameter. This intuition is confirmed in Figure 11a. Using \hat{T}_k^0 yields significantly more ambiguous results than an equivalent procedure based on \hat{T}_k^2 . Decreasing the threshold for the eigenvalues leads to a gradual increase in the number of estimated factors, indicated by the absence of a prolonged flat region. Using a thresholding rule based on \hat{T}_k^2 instead, we observe two flat regions in the graph at $\hat{r} = 3$ and $\hat{r} = 6$ respectively.

Alternatively, a practitioner could vary τ and observe how the estimated number of factors changes. The parameter τ can here be thought of as a complexity parameter because varying τ shifts the steep region of Figure 3. I demonstrate this in Figure 11b, which depicts the value of the TR estimator as a function of the complexity parameter τ . Figure 11b suggests that there is one “global” factor in the data and that the second most important factor is significantly weaker than the first one, as indicated by the first flat region of the graph. Next, the 7th factor appears to be significantly weaker than the 6th, as indicated by the second flat region in Figure 11b. Thus, Figure 11 suggests that the choice of $\hat{r} = 6$ is quite robust when incorporating the additional information



(a) 25 largest loadings for factors 2-7. In line with Figure 10b, I see that factors 3, 5 and 6 are more concentrated on a subset of outcomes.



(b) Incremental R^2 of the common component for each series, ordered from largest to smallest, for adding factors 5 and 6 to the model (blue line). For comparison the dotted red line depicts adding the two next factors

Figure 12: Both figures illustrate how factors 5 and 6 are highly locally concentrated.

in the eigenvectors.

My findings also suggest a more nuanced interpretation of the factors: While factors two and four appear to weakly affect a large cross section of outcomes, I find evidence that factors 3, 5 and 6 are more “local” in the sense that their loadings are concentrated on a subset of outcomes. Figure 12a illustrates this. For each factor, the associated loadings were ordered (in absolute value). Figure 12a then depicts how the largest 25 loadings decay for factors 3-8. The three factors I identified as being “local” can easily be distinguished. They exhibit some large loadings, combined with a steeper subsequent drop-off.

I next illustrate the importance of individual factors by considering the R^2 of the common component in explaining movements in various series. For a given series in the panel, this statistic measures the variation in the series due to contemporaneous variation in the factors.¹⁸ The addition of factors 5 and 6 to the model increases the R^2 of the common component for each series on average by 0.06. The smooth nature of the screeplot implies that the estimated 7th and 8th factors are not significantly weaker and adding them to the model produces an incremental R^2 of the common component for each series of 0.05 on average, barely less than the increase that resulted from adding the fifth and sixth factor. Figure 12b depicts how the addition of weaker factors affects the individual R^2 of the 94 series in the panel. Sorted from highest to lowest, it illustrates how the local factors 5 and 6 are particularly important for some of the outcomes. For example, for the

¹⁸Another common performance measure is the h-step ahead forecasting performance for a number of either composite indices or individual series from the panel. I pursue this avenue in a separate paper that focuses on the estimation of the factor space (Freyaldenhoven (2017)).

Factor 3	Factor 5	Factor 6
PPI: Int. Material: Supplies & Components PPI: Industrial Commodities PPI: Finished Consumer Goods PPI: Crude Petroleum Defl by PCE(LFE) Gasoline and other energy goods BS: Real Compensation Per Hour Nonfarm: Real Compensation Per Hour BS: Implicit Price Deflator ISM Manufacturing: Prices Paid Index	Nonfarm: Unit Nonlabor Payments Nonfarm: Unit Labor Cost Nonfarm: Real Compensation Per Hour BS: Real Compensation Per Hour PPI: Finished Consumer Foods Food & beverages for off-premises consump Nonfarm: Output Per Hour of All Persons PPI: Finished Consumer Goods	tb6m-tb3m GS1-Tb3m GS10-Tb3m S&P'S STOCK PRICE INDEX DOW JONES IA Consumer Loans, All Commercial Banks BAA-GS10 Spread

Table 4: Variables corresponding to largest loadings for Factors 3,5 and 6, the most local factors according to Figure 10. Red coloring indicates negative loading, while black refers to positive loadings.

most impacted series of the panel, factors 5 and 6 explain around 70 per cent of the variation in that series. Figure 12b therefore demonstrates that a subset of the series is very well explained by factors 5 and 6, associated with a large jump in the corresponding R^2 , whereas the addition of factors 7 and 8 has no such strong effect on any individual series. The 10 largest increments in the model R^2 are on average 0.39 and 0.18 for factors 5+6 and 7+8 respectively. Unsurprisingly the gains are most extreme for the series associated with the largest loadings on factors 5 and 6.

A further appealing implication of my framework is that, by treating factors as local the resulting factors may be easier to interpret as they only correspond to a small subset of the observables, contrasting with conventional factors, which are often hard to interpret. Table 4 shows which economic indicators correspond to the largest loadings (in absolute value) associated with the three local factors. Variables with a negative loading are shown in red. For the first column (factor 3), I note that 6 of the 9 variables, printed in bold, represent price indices as classified in the handbook chapter of Stock and Watson (2016). Additionally the fourth entry, while classified separately as an ‘Oil market variable’, also represents a price index. The remaining two variables are both classified as ‘Productivity and Earnings’ and it is worth noting that they have the opposite sign. Among the 94 indicators considered in this exercise are 5 that are classified as ‘Productivity and Earnings’ in Stock and Watson (2016). All 5 of these appear in the second column of table 4, emphasized in bold. Further, the remaining 3 entries are all price indices. I also note that the 6th factor is highly concentrated on a number of financial variables, specifically spreads and stock market indicators (again emphasized in bold). Further, this factor is associated with a negative return on the stock market and an increase in the interest rate spread.

This aids in the interpretation of the factors: For example, based on the discussion above, the 6th factor could be interpreted as indicating a flight from stocks into safe assets, such as bonds.

The previous discussion illustrates the advantage of taking the eigenvectors into account when selecting the number of factors as proposed in this paper. Without this additional information, factors 5 and 6 are missed by two of the three existing estimators in the literature. But these factors are highly influential on a subset of outcomes as shown above. Failure to include them in the model would thus result in a model that does very poorly in explaining this part of the economy.

7 Concluding Remarks

In this paper I develop a framework for factor models that allows for local factors which only affect an unknown subset of the observables. In many economic models I find that factors affecting proportionally more than \sqrt{n} of the n observed variables are of economic interest. Under standard assumptions on the error terms, this coincides with the number of factors can be estimated consistently using the principal component estimator. I further show that existing estimators for the number of factors in general do not yield a consistent estimate for this number of “relevant” factors. To estimate the number of economically important and estimable factors consistently I argue that there is additional information in the eigenvectors that has not been exploited in the past. I demonstrate how one can incorporate this information into some of the prominent estimators commonly used. Monte Carlo evidence suggests significant finite sample gains over those existing estimators.

In cases in which there is no clear gap in the distribution of eigenvalues, the theory developed in this paper provides a viable framework. It further provides a theoretical foundation that justifies the use of both factor models and the principal component estimator in datasets with no such clear gap.

In addition, the methods of this paper provide a novel insight into the structure of the data. There are two potential reasons subsequent factors may appear “weak” in a given dataset - either a weaker factor can have a weak effect on all observables, or it can have a strong impact on only a subset of observables (which I call a “local” factor in this paper). By using only the eigenvalues of $\frac{X'X}{T}$, these two kinds of factors will be treated equally. However, depending on the economic model or context, a researcher may be more interested in one or the other. By incorporating information from the eigenvectors, I allow a practitioner to distinguish between the two cases.

Finally, I implement my methods in one of the canonical datasets used in the factor model literature and find strong evidence that there are indeed local factors present in the data.

The analysis in this paper suggests a number of promising topics for future research. Perhaps most interestingly, I conjecture that the principal component estimator considered in this paper can be substantially improved upon (at least in finite samples) using the sparsity assumptions of the model for the estimation of the factors. A regularized estimation approach suggests itself and is currently investigated in a separate project (Freyaldenhoven (2017)).

A Monte Carlo Simulation

A.1 DGP for table 1

$$\underset{(500 \times n)}{X} = \underset{(500 \times 2)(2 \times n)}{F} \underset{(2 \times n)}{\Lambda^T} + \underset{(500 \times n)}{e}. \quad (86)$$

I observe a panel with $T = 500$, where the cross-sectional dimension varies across simulations (see table 1). The variables exhibit the following factor structure: 2 independent factors $F_k, k = 1, 2$ are created from a standard normal distribution with an autocorrelation of .3. Λ is a matrix of ones and zeros such that:

$$\begin{aligned} X_j &= F_1 + e_j, & \text{for } j \in \mathcal{A}_2^c \\ X_j &= F_2 + e_j, & \text{for } j \in \mathcal{A}_2 \end{aligned}$$

The cardinality of \mathcal{A}_2 is varied from $n^{\frac{1}{4}}$ to $n^{\frac{3}{4}}$ as depicted in table 1.

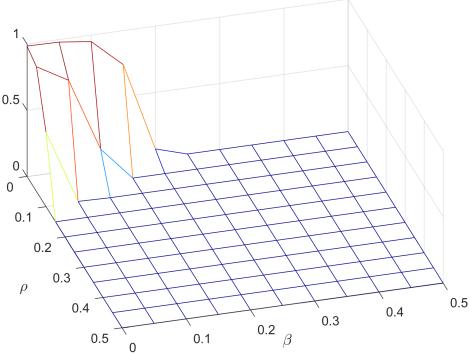
Finally, I allow the idiosyncratic errors to exhibit both cross-sectional as well as inter-temporal correlation. I follow Onatski (2010), Onatski (2012) and model the errors as follows:

$$e_{ti} = \rho e_{t-1,i} + (1 - \rho^2)^{1/2} v_{it} \quad (87)$$

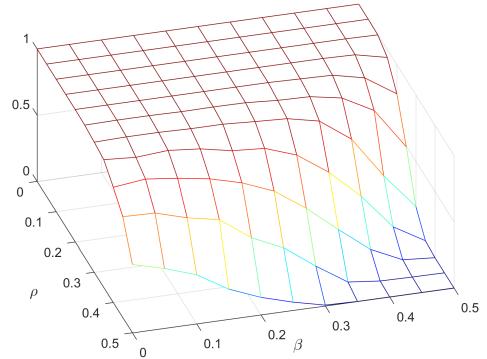
$$v_{ti} = \beta v_{t,i-1} + (1 - \beta^2)^{1/2} u_{it}, \quad u_{it} \sim i.i.d N(0, 1). \quad (88)$$

I set both ρ and β to .3 to allow for modest correlations in the error terms.

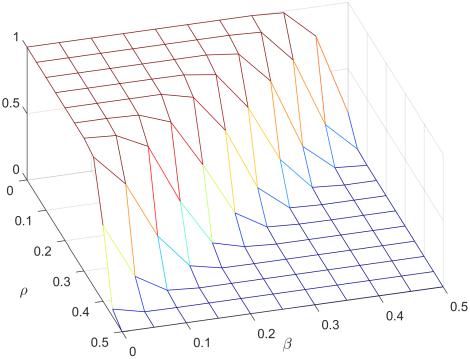
A.2 Results under unfavorable DGP



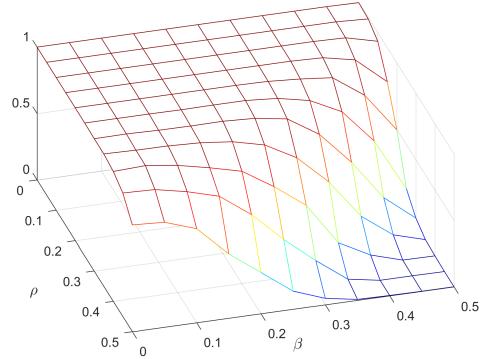
(a) Information criterion PC from Bai and Ng (2002) (PC)



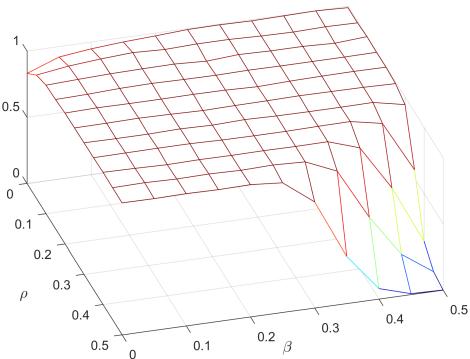
(b) Thresholding based on T^2 (TC)



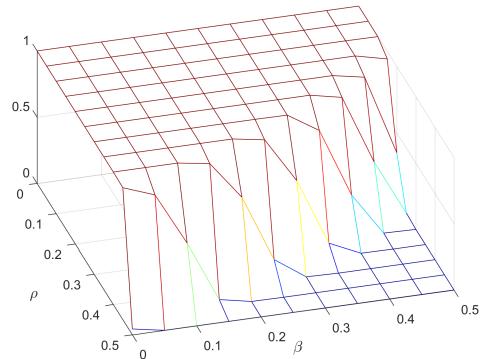
(c) Maximum ratio of two subsequent eigenvalues as in Ahn and Horenstein (2013) (ER)



(d) Maximum ratio of two subsequent values of T^2 (TR)

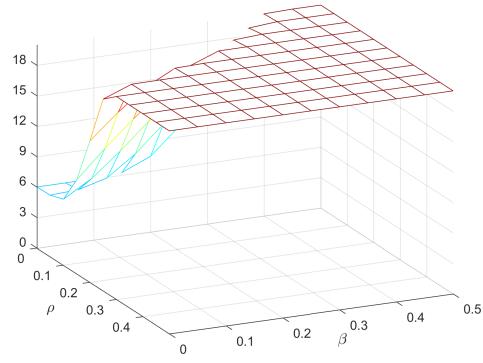


(e) Difference of two subsequent eigenvalues as in Onatski (2010) (ED)

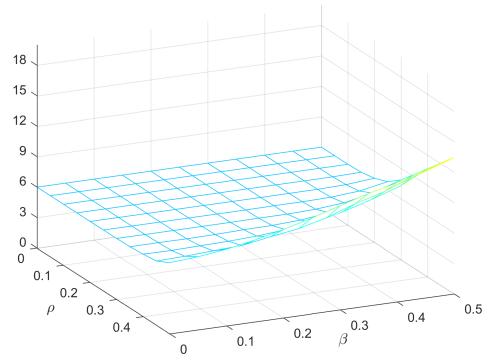


(f) Thresholding based on T^0 ($PC_{\sqrt{n}}$)

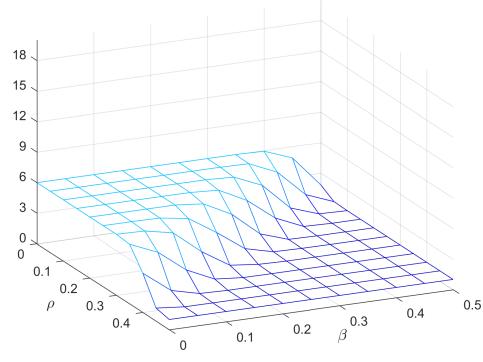
Figure 13: Percentage of correctly estimating the true number of factors as both cross sectional and intertemporal correlation is varied. $r = 6$, $n = 300$, $T = 500$, $(\rho, \beta) = (0.3, 0.1)$, $\theta = 1.5$. For each entry in Λ , $\lambda_{ik} = 1 + \nu_{ik}$, where $\eta_{ik} \sim N(0, 1)$. Figure based on 500 replications.



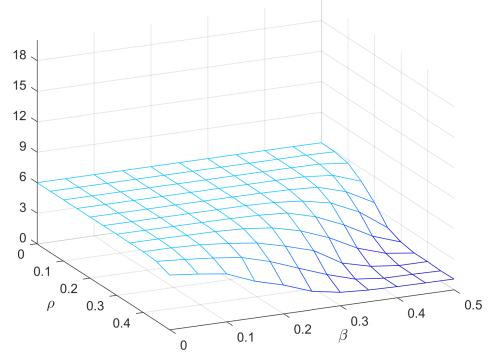
(a) PC criterion from Bai and Ng (2002) (PC)



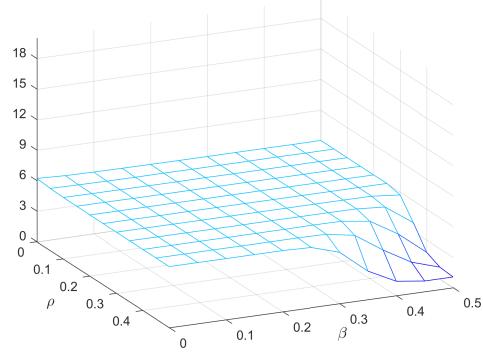
(b) Thresholding based on T^2 (TC)



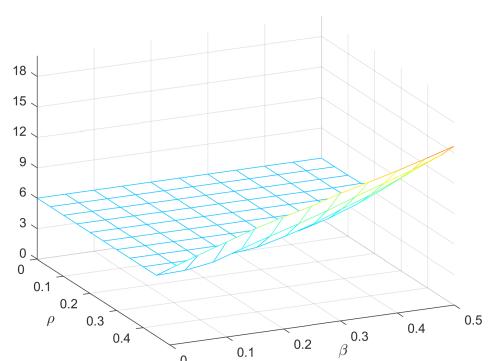
(c) Maximum ratio of two subsequent eigenvalues as in Ahn and Horenstein (2013) (ER)



(d) Maximum ratio of two subsequent values of T^2 (TR)

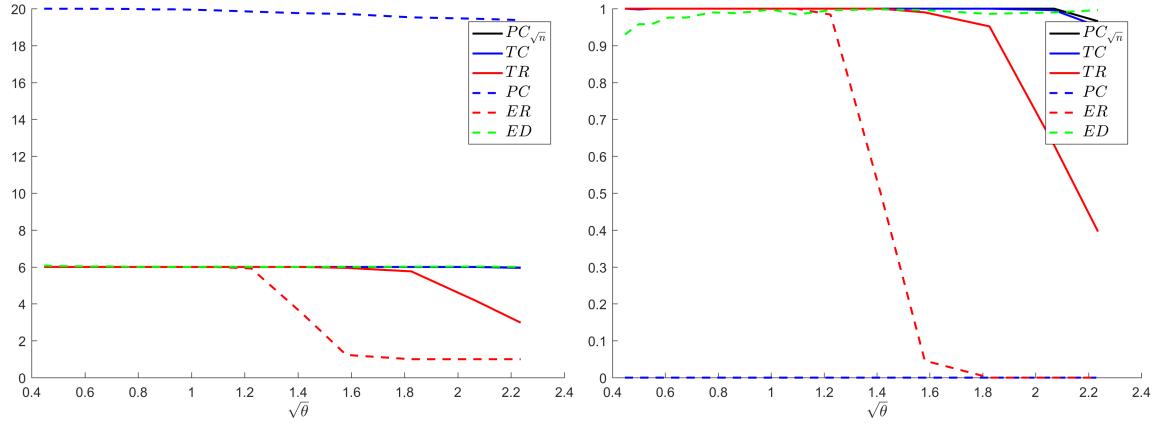


(e) Difference of two subsequent eigenvalues as in Onatski (2010) (ED)



(f) Thresholding based on T^0 ($PC_{\sqrt{n}}$)

Figure 14: Average number of factors estimated as both cross sectional and intertemporal correlation is varied. $r = 6$, $n = 300$, $T = 500$, $(\rho, \beta) = (0.3, 0.1)$, $\theta = 1.5$. For each entry in Λ , $\lambda_{ik} = 1 + \nu_{ik}$, where $\eta_{ik} \sim N(0, 1)$. Figure based on 500 replications.



(a) Average estimated number of factors according to TC , TR , PC , ER and ED (b) % of correctly estimated number of factors according to TC , TR , PC , ER and ED

Figure 15: Empirical behavior of estimators as the relative variance of idiosyncratic noise increases. $r = 6$, $n = 300$, $T = 500$, $(\rho, \beta) = (0.3, 0.1)$, $\theta = 1.5$. For each entry in Λ , $\lambda_{ik} = 1 + \nu_{ik}$, where $\eta_{ik} \sim N(0, 1)$. Figure based on 500 realizations.

<i>n</i>	<i>T</i>	<i>ER</i>	<i>TR</i>	<i>PC</i>	<i>PC</i> _{\sqrt{n}}	<i>TC</i>	<i>ED</i>
100	100	1 / 0.00	2.89 / 0.27	20 / 0.00	15.9 / 0.00	11.8 / 0.00	4.7 / 0.69
100	150	1 / 0.00	3.35 / 0.45	20 / 0.00	13 / 0.00	9.61 / 0.00	5.82 / 0.94
150	100	1 / 0.00	3.52 / 0.46	20 / 0.00	15.1 / 0.00	11.9 / 0.00	5.64 / 0.90
150	250	1.16 / 0.03	5.62 / 0.92	20 / 0.00	7.13 / 0.14	6.54 / 0.60	6.02 / 0.99
150	500	2.3 / 0.26	5.89 / 0.97	19.5 / 0.00	6 / 1.00	6.03 / 0.97	6 / 1.00
300	250	2.32 / 0.26	5.99 / 1.00	20 / 0.00	6 / 1.00	6.01 / 0.99	6 / 1.00
300	500	5.85 / 0.97	6 / 1.00	19.9 / 0.00	6 / 1.00	6 / 1.00	6.01 / 0.99
300	750	5.99 / 1.00	6 / 1.00	17.6 / 0.00	6 / 1.00	6 / 1.00	6.01 / 0.99
500	250	4.28 / 0.66	6 / 1.00	20 / 0.00	6 / 1.00	6 / 1.00	6 / 1.00
500	500	6 / 1.00	6 / 1.00	20 / 0.00	6 / 1.00	6 / 1.00	6 / 1.00
500	750	6 / 1.00	6 / 1.00	18.5 / 0.00	6 / 1.00	6 / 1.00	6 / 1.00
1000	1000	6 / 1.00	6 / 1.00	17.6 / 0.00	6 / 1.00	6 / 1.00	6.01 / 0.99

Table 5: Each entry depicts a combination $\hat{r}/\%$, where \hat{r} is the average number of estimated factors, and $\%$ is the percentage correctly classifying $r = 6$ based on 500 replications. For each entry in Λ , $\lambda_{ik} = 1 + \nu_{ik}$, where $\eta_{ik} \sim N(0, 1)$. In each row, the highest percentage is highlighted.

B Auxiliary Figures

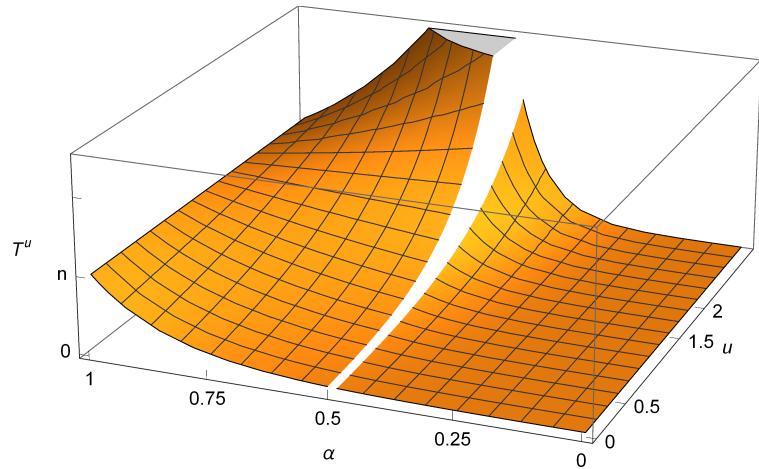


Figure 16: Theoretical divergence rate of $T_{z_k}^u$ as a function of both factor strength α and tuning parameter u extending past $u = 2$ for $\tau = .5$.

For example, with $u = 3$:

$$T_{z_k}^u = \begin{cases} O_p(n^{\frac{3}{2} - \frac{1}{2}\alpha_k}) & \text{for } \alpha_k \geq .5 \\ O_p(n^{\frac{5}{2}\alpha_k})g(n)^{-3} & \text{otherwise} \end{cases}$$

C The Effect of Including Aggregates on Factor Analysis

To see why the inclusion of series at multiple levels of aggregation when estimating a factor model is problematic, consider the following simple example of two time series following a (strict) factor structure as follows:

$$\begin{aligned}X_1 &= F + e_1 \\X_2 &= F + e_2 \quad , e_1 \perp e_2.\end{aligned}$$

A researcher wishing to estimate a factor model on her data includes a third variable $X_3 = (X_1 + X_2)/2$, a simple composite index of the two first variable with equal weights. Clearly $X_3 = (X_1 + X_2)/2 = F + (e_1 + e_2)/2$. Thus, a factor representation of the larger set would look as follows:

$$\begin{aligned}X_1 &= F + e_1 \\X_2 &= F + e_2 \\X_3 &= F + e_3.\end{aligned}$$

However, for factor analysis to be successful using the principal component estimator, the key is no (or low) correlation in the errors. Clearly, this is no longer the case, since $e_3 = (e_1 + e_2)/2$ by construction.

D Mathematical Appendix

D.1 Auxiliary Lemmata

Lemma 3. Under Assumptions 1-3, for all n and T

$$\frac{1}{T} \sum_{s=1}^T \sum_{t=1}^T \mathbb{E} \left(\frac{e'_s e_t}{n} \right)^2 \leq C \quad (89)$$

Proof. See Lemma 1(i) in Bai and Ng (2002), using Assumption 3 (b). \square

Lemma 4. Under Assumptions 1-4, for any fixed K , let A be a $T \times K$ matrix A such that $A'A = T I_K$. Define $\alpha_{max} = \max_k \alpha_k$, $k = 1, 2, \dots, K$.

Then:

$$\frac{1}{T^2} |trace(A' F \Lambda' e' A)| = O_p(n^{\frac{1}{2}\alpha_{max}}) \quad (90)$$

Proof.

$$|trace(A' F \Lambda' e' A)| = \left| trace(A' [\sum_{k=1}^r F_k \lambda'_{.k} e'] A) \right| \quad (91)$$

$$= \left| \sum_{k=1}^r trace(A' [F_k \lambda'_{.k} e'] A) \right| \quad (92)$$

$$\leq \left| \sum_{k=1}^r \|A'\| \|F_k\| \|\lambda'_{.k} e'\| \|A\| \right| \quad (93)$$

$$= \left| \sum_{k=1}^r \|A\|^2 \|F_k\| \sqrt{\sum_t (\sum_i \lambda_{ik} e_{it})^2} \right| \quad (94)$$

By Assumption 4(a) the inner most sum grows at an order of $O_p(n^{\frac{1}{2}\alpha_k})$. I conclude:

$$\frac{1}{T^2} |trace(A' F \Lambda' e' A)| \leq \left| \sum_{k=1}^r \left\| \frac{1}{\sqrt{T}} A \right\|^2 \left\| \frac{1}{\sqrt{T}} F_k \right\| \sqrt{\frac{1}{T} \sum_t (\sum_i \lambda_{ik} e_{it})^2} \right| \quad (95)$$

$$= \left| \sum_{k=1}^r O_p(1) \times O_p(n^{\frac{1}{2}\alpha_k}) \right| \quad (96)$$

$$= O_p(n^{\frac{1}{2}\alpha_{max}}), \quad (97)$$

which completes the proof. I note that in most cases at least one factor will be strong, corresponding to $\alpha_{max} = 1$. In that case, the above rate becomes $O_p(\sqrt{n})$.

□

Lemma 5. Under assumptions 1- 4 , for any fixed K , let A be a $T \times K$ matrix A such that $A'A = TI_K$. Define $\alpha_{max} = \max_k \alpha_k$, $k=1, 2, \dots, K$.

$$\sup_A (A' \frac{XX'}{T} A - A' \frac{F\Lambda'\Lambda F'}{T} A) = O_p(n^{\frac{1}{2}\alpha_{max}}) \quad (98)$$

Proof.

$$\sup_A (A' \frac{XX'}{T^2} A - A' \frac{F\Lambda'\Lambda F'}{T^2} A) = \sup_A A' \left(\frac{e\Lambda F'}{T^2} + \frac{F\Lambda'e'}{T^2} + \frac{ee'}{T^2} \right) A \quad (99)$$

$$\leq \sup_A A' \left(\frac{e\Lambda F'}{T^2} + \frac{F\Lambda'e'}{T^2} \right) A + \sup_A A' \frac{ee'}{T^2} A \quad (100)$$

$$= O_p(n^{\frac{1}{2}\alpha_{max}}) + O_p(1) \quad (101)$$

where the last equality follows from Lemma 4 and Assumption 3(e)

□

Lemma 6. Let $\hat{F}_1 = \arg \sup_\gamma \gamma' \frac{XX'}{T} \gamma$, such that $\frac{\gamma' \gamma}{T} = 1$. In other words: Let \hat{F} be the standardized eigenvector corresponding to the largest eigenvalue of XX' . Then, under Assumptions 1- 4:

$$\frac{\hat{F}'_1 F_1}{T} = 1 + O_p(n^{-\frac{1}{2}\alpha_{max}}), \quad (102)$$

and for $k = 2, \dots, r$

$$\frac{\hat{F}'_1 F_k}{T} = O_p(n^{-\frac{1}{4}\alpha_{max}}), \quad (103)$$

Proof. Decompose \hat{F}_1 as follows:

$$\hat{F}_1 = F \left(\frac{F' F}{T} \right)^{-\frac{1}{2}} \xi_1 + \hat{V} \text{ such that } \hat{V}' F = 0 \quad (104)$$

Since $\frac{\hat{F}'_1 \hat{F}_1}{T} = \xi'_1 \xi_1 + \frac{\hat{V}' \hat{V}}{T}$, this implies $\xi'_1 \xi_1 \leq 1$. Then:

$$\hat{F}'_1 \frac{F \Lambda' \Lambda F'}{T^2} \hat{F}_1 = [F(\frac{F' F}{T})^{-\frac{1}{2}} \xi_1 + \hat{V}]' \left(\frac{F \Lambda' \Lambda F'}{T^2} \right) [F(\frac{F' F}{T})^{-\frac{1}{2}} \xi_1 + \hat{V}] \quad (105)$$

$$= \xi'_1 \left(\frac{F' F}{T} \right)^{\frac{1}{2}} \Lambda' \Lambda \left(\frac{F' F}{T} \right)^{\frac{1}{2}} \xi_1 \quad (106)$$

$$= \xi'_1 I_r D_r^{(n)} I_r \xi_1 \quad (107)$$

by assumptions 1(b) and (c). I also note that

$$\begin{aligned} \frac{1}{T^2} (\hat{F}'_1 F \Lambda' \Lambda F' \hat{F}_1 - \hat{F}'_1 X X' \hat{F}_1) &= \frac{1}{T^2} (\hat{F}'_1 F \Lambda' \Lambda F' \hat{F}_1 - F'_1 F \Lambda' \Lambda F' F_1) \\ &\quad + \frac{1}{T^2} (F'_1 F \Lambda' \Lambda F' F_1 - \hat{F}'_1 X X' \hat{F}_1) \end{aligned} \quad (108)$$

$$= O_p(n^{\frac{1}{2}\alpha_{max}}) \quad (\text{by Lemma 5}) \quad (109)$$

The second term on the RHS is simply the difference between the largest eigenvalue of $X X' / T$ and $F \Lambda' \Lambda F' / T$. By Theorem 1 that difference is $O_p(n^{\frac{1}{2}\alpha_{max}})$. It follows that the first term on the RHS is also $O_p(n^{\frac{1}{2}\alpha_{max}})$. I therefore obtain

$$\frac{1}{T^2} (\hat{F}'_1 F \Lambda' \Lambda F' \hat{F}_1 - F'_1 F \Lambda' \Lambda F' F_1) = \xi'_1 D_r^{(n)} \xi_1 - d_1 \quad (110)$$

$$= (\xi_{11}^2 - 1) d_1 + \sum_{k=2}^r \xi_{1k}^2 d_k \quad (111)$$

$$= O_p(n^{\frac{1}{2}\alpha_{max}}) \quad (112)$$

Because $d_1 > d_2 > \dots > 0$, this implies that $\xi_{11}^2 - 1 = O_p(n^{-\frac{1}{2}\alpha_{max}})$. Since $\xi' \xi \leq 1$ it follows that $\xi_{1k}^2 = O_p(n^{-\frac{1}{2}\alpha_{max}})$ for $k = 2, \dots, r$ and $\frac{\hat{V}' \hat{V}}{T} = O_p(n^{-\frac{1}{2}\alpha_{max}})$. Since also (from the decomposition in 104)

$$\frac{\hat{F}'_1 F}{T} = \frac{1}{T} [F(\frac{F' F}{T})^{-1/2} \xi_1 + \hat{V}]' F \quad (113)$$

$$= \xi'_1 \left(\frac{F' F}{T} \right)^{-1/2} \frac{F' F}{T} \quad (114)$$

$$= \xi'_1, \quad (115)$$

$\left(\frac{\hat{F}'_1 F_1}{T} \right)^2 - 1 = O_p(n^{-\frac{1}{2}\alpha_{max}})$ and therefore $\frac{\hat{F}'_1 F_1}{T} = 1 + O_p(n^{-\frac{1}{2}\alpha_{max}})$. It then also follows that, for $k = 2, \dots, r$: $\frac{\hat{F}'_1 F_k}{T} = O_p(n^{-\frac{1}{4}\alpha_{max}})$.¹⁹ \square

¹⁹Here, and throughout the paper, I assume that the sign indeterminacy of \hat{F} is resolved by adding the normalization

Lemma 7. Let \hat{F} be the $T \times K$ matrix of the first K eigenvectors of XX' , normalised such that $\frac{\hat{F}'\hat{F}}{T} = I_K$. Then, under Assumptions 1-4, for each $k = 1, \dots, K$ and $l = 1, \dots, r$:

- For $k < l$: $\frac{\hat{F}'_k F_l}{T} = \bar{O}_p(n^{\frac{1}{4}\alpha_{max} - \frac{1}{2}\alpha_k})$
- For $k = l$: $\frac{\hat{F}'_k F_l}{T} = 1 + \bar{O}_p(n^{\frac{1}{2}\alpha_{max} - \alpha_l})$
- For $k > l$: $\frac{\hat{F}'_k F_l}{T} = \bar{O}_p(n^{\frac{1}{4}\alpha_{max} - \frac{1}{2}\alpha_l})$

Proof. The result for the first row of $\frac{\hat{F}'F}{T}$ is given in Lemma 6. For the remaining columns we repeat the steps above in orthonormal subspaces. My strategy is therefore similar to the one followed in Stock and Watson (2002a). However, allowing for varying factor strengths requires a more nuanced consideration of the subsequent principal components. Additionally, unlike Stock and Watson (2002a), I explicitly derive the rates of convergence for all quantities of interest.

Using the same reasoning as in the previous lemma , we obtain a representation of \hat{F}_k , the k th column of \hat{F} , as follows:

$$\hat{F}_k = F\left(\frac{F'F}{T}\right)^{-\frac{1}{2}}\xi_k + \hat{V}_k \text{ such that } \hat{V}'_k F = \mathbf{0}. \quad (116)$$

This implies $\xi'_k \xi_k \leq 1$,

$$\hat{F}'_k \frac{F \Lambda' \Lambda F'}{T^2} \hat{F}_k = \xi'_k I_r D_r I_r \xi_k \quad (117)$$

and

$$\frac{1}{T^2} (\hat{F}'_k F \Lambda' \Lambda F' \hat{F}_k - \hat{F}'_k X X' \hat{F}_k) \quad (118)$$

$$= \frac{1}{T^2} (\hat{F}'_k F \Lambda' \Lambda F' \hat{F}_k - F'_k F \Lambda' \Lambda F' F_k) + \frac{1}{T^2} (F'_k F \Lambda' \Lambda F' F_k - \hat{F}'_k X X' \hat{F}_k) \quad (119)$$

$$= O_p(n^{\frac{1}{2}\alpha_{max}}). \quad (120)$$

By Theorem 1 the second term on the RHS is $O_p(n^{\frac{1}{2}\alpha_k})$. For the k th largest eigenvector this implies for the first term

$$\frac{1}{T^2} (\hat{F}'_k F \Lambda' \Lambda F' \hat{F}_k - F'_k F \Lambda' \Lambda F' F_k) = \xi'_k D_r \xi_k - d_k \quad (121)$$

$$= (\xi_{kk}^2 - 1)d_k + \sum_{l \neq k} \xi_{kl}^2 d_l \quad (122)$$

$$= O_p(n^{\frac{1}{2}\alpha_{max}}) \quad (123)$$

that the diagonal elements of $\hat{F}'F$ are nonnegative.

Because $d_1 > d_2 > \dots > 0$, this implies that

- For $k = l$: $\xi_{kk}^2 = 1 + O_p(n^{\frac{1}{2}\alpha_{max} - \alpha_l})$
- For $k < l$: $\xi_{kl}^2 = O_p(n^{\frac{1}{2}\alpha_{max} - \alpha_k})$
- For $k > l$: $\xi_{kl}^2 = O_p(n^{\frac{1}{2}\alpha_{max} - \alpha_l})$ (Since $\xi'_k \xi_k \leq 1$).

We further note that $\frac{\hat{F}'_k F}{T} \leq 1$ and hence $\xi_{kl}^2 = O_p(1) \forall l$. This also implies a lower bound on the factor strength, indicated by α_k , for which ξ_k^2 is guaranteed to converge: $\alpha_k > \frac{1}{2}\alpha_{max}$. \square

Lemma 8. Let \hat{F} be the $T \times K$ matrix of the first K eigenvectors of XX' , normalised such that $\frac{\hat{F}' \hat{F}}{T} = I_K$ and define a $(r \times K)$ matrix $H = \Lambda' \Lambda \frac{F' \hat{F}}{T} \hat{D}_K^{-1}$, where \hat{D}_K is a diagonal matrix with the K largest eigenvalues of $\frac{X' X}{T}$ on the main diagonal. Then, under Assumptions 1-4:

$$\frac{1}{T} \sum_{t=1}^T (\hat{F}_{tk} - H'_{k \cdot} F_t)^2 = \bar{O}_p(n^{1-2\alpha_k}), \quad (124)$$

where $H'_{k \cdot}$ denotes the k 'th row of H' .

Proof. Note that by the properties of eigenvectors and eigenvalues $\hat{F} = \frac{XX'}{T} \hat{F} \hat{D}_K^{-1}$. Then:

$$\hat{F} - FH = \frac{XX'}{T} \hat{F} \hat{D}_K^{-1} - F \Lambda' \Lambda \frac{F' \hat{F}}{T} \hat{D}_K^{-1} \quad (125)$$

$$= \frac{1}{T} (XX' - F \Lambda' \Lambda F') \hat{F} \hat{D}_K^{-1} \quad (126)$$

$$= \frac{1}{T} (ee' + e \Lambda F' + F \Lambda' e') \hat{F} \hat{D}_K^{-1} \quad (127)$$

This is related to the decomposition first derived in Bai and Ng (2002) and used extensively in the literature since its introduction (e.g. Bai (2003), Choi (2012)). The following derivations therefore follow those in Bai and Ng (2002) and Bai (2003), who consider only strong factors. For a particular t we may write:

$$\hat{F}_t - H' F_t = \frac{1}{T} \hat{D}_K^{-1} \hat{F}' (ee_t + e \Lambda F_t + F \Lambda' e_t) \quad (128)$$

$$= \hat{D}_K^{-1} \left(\frac{1}{T} \sum_{s=1}^T \hat{F}_s e'_s e_t + \frac{1}{T} \sum_{s=1}^T \hat{F}_s F'_s \Lambda' e_t + \frac{1}{T} \sum_{s=1}^T \hat{F}_s e'_s \Lambda F_t \right) \quad (129)$$

Because $(I + II + III)^2 \leq 3(I^2 + II^2 + III^2)$, by Cauchy-Schwarz and submultiplicity of the

norm: $\|\hat{F}_t - H'F_t\|^2 \leq \|\hat{D}_r^{-1}\|^2 3(I_t + II_t + III_t)$, where:

$$I_t = \frac{1}{T^2} \left\| \sum_{s=1}^T \hat{F}_s e'_s e_t \right\|^2 \quad (130)$$

$$II_t = \frac{1}{T^2} \left\| \sum_{s=1}^T \hat{F}_s F'_s \Lambda' e_t \right\|^2 \quad (131)$$

$$III_t = \frac{1}{T^2} \left\| \sum_{s=1}^T \hat{F}_s e'_s \Lambda F_t \right\|^2 \quad (132)$$

Thus $\frac{1}{T} \sum_{t=1}^T \|\hat{F}_t - H'F_t\|^2 \leq \|\hat{D}_K^{-1}\|^2 \frac{1}{T} \sum_{t=1}^T 3(I_t + II_t + III_t)$, while for each individual factor estimate \hat{F}_k , $k = 1, 2, \dots, r$, $\frac{1}{T} \sum_{t=1}^T (\hat{F}_{tk} - H'_k F_{tk})^2 \leq \|\hat{d}_k^{-1}\|^2 \frac{1}{T} \sum_{t=1}^T 3(I_{kt} + II_{kt} + III_{kt})$, with the r -by-1 vector \hat{F}_s replaced by the scalar \hat{F}_{sk} in each of I_t , II_t and III_t above.

Consider each of the above three terms separately:

$$\frac{1}{T} \sum_{t=1}^T I_{tk} = \frac{1}{T} \sum_{t=1}^T \left\| \frac{1}{T} \sum_{s=1}^T \hat{F}_{sk} e'_s e_t \right\|^2 \quad (133)$$

$$\leq \frac{1}{T} \sum_{t=1}^T \left(\left\| \frac{1}{T} \sum_{s=1}^T \hat{F}_{sk} [e'_s e_t - \mathbb{E}(e'_s e_t)] \right\|^2 + \left\| \frac{1}{T} \sum_{s=1}^T \hat{F}_{sk} \mathbb{E}(e'_s e_t) \right\|^2 \right) \quad (134)$$

$$= O_p(n) \quad (135)$$

Since this part does not involve any non-standard assumptions (it does not involve the factor loadings), the last equality follows using the same arguments as in the proof of Theorem 1 in Bai and Ng (2002) using Lemma 3 and Assumption 3(c). Details are not worth repeating. For the next part:

$$\frac{1}{T} \sum_{t=1}^T II_{tk} = \frac{1}{T} \sum_{t=1}^T \frac{1}{T^2} \left\| \sum_{s=1}^T \hat{F}_{sk} F'_s \Lambda' e_t \right\|^2 \quad (136)$$

$$\leq \frac{1}{T} \sum_{t=1}^T \left[\|\Lambda' e_t\|^2 \left(\frac{1}{T} \sum_{s=1}^T \|F_s\|^2 \right) \left(\frac{1}{T} \sum_{s=1}^T \|\hat{F}_{sk}\|^2 \right) \right] \quad (137)$$

$$= \frac{1}{T} \sum_{t=1}^T \left[\|\Lambda' e_t\|^2 O_p(1) \right] \quad (138)$$

$$= O_p(n^{\alpha_{max}}), \quad (139)$$

by Assumption 4(a). Finally, for III_{kt} one can show in a similar manner that

$$\frac{1}{T} \sum_{t=1}^T III_{kt} = \frac{1}{T} \sum_{t=1}^T \frac{1}{T^2} \left\| \sum_{s=1}^T \hat{F}_{sk} e'_s \Lambda F_t \right\|^2 \quad (140)$$

$$= O_p(n^{\alpha_{max}}) \quad (141)$$

Consequently

$$\frac{1}{T} \sum_{t=1}^T \left\| \hat{F}_{tk} - H'_{k \cdot} F_t \right\|^2 \leq \hat{d}_k^{-2} 3(I_{kt} + II_{kt} + III_{kt}) \quad (142)$$

$$\leq O_p(n^{-2\alpha_k}) \left(O_p(n) + O_p(n^{\alpha_{max}}) + O_p(n^{\alpha_{max}}) \right) \quad (143)$$

$$= O_p(n^{1-2\alpha_k}) \quad (144)$$

□

Lemma 9. Define a matrix $H = \Lambda' \Lambda \frac{F' \hat{F}}{T} \hat{D}_K^{-1}$, where \hat{D}_K is a diagonal matrix with the K largest eigenvalues of $\frac{X' X}{T}$ on the main diagonal. Accordingly, let $H_{\cdot k}$ denote the k th column of H . Then, under Assumptions 1-4, $H_{\cdot k} = \iota_k + O_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2}-\alpha_k})$.

Proof. First note that

$$H = \Lambda' \Lambda \frac{F' \hat{F}}{T} \hat{D}_K^{-1} = \begin{bmatrix} \frac{d_1}{\hat{d}_1} \frac{F'_1 \hat{F}_1}{T} & \frac{d_1}{\hat{d}_2} \frac{F'_1 \hat{F}_2}{T} & \dots & \frac{d_1}{\hat{d}_K} \frac{F'_1 \hat{F}_K}{T} \\ \frac{d_2}{\hat{d}_1} \frac{F'_2 \hat{F}_1}{T} & \frac{d_2}{\hat{d}_2} \frac{F'_2 \hat{F}_2}{T} & & \vdots \\ \vdots & & \ddots & \\ \frac{d_r}{\hat{d}_1} \frac{F'_r \hat{F}_1}{T} & \dots & \frac{d_r}{\hat{d}_K} \frac{F'_r \hat{F}_K}{T} & \end{bmatrix}, \quad (145)$$

where d_k and \hat{d}_k denote the k th entry on the diagonal of $\Lambda' \Lambda$ and \hat{D}_K respectively. Consider entry H_{lk} at position (l, k) . First note that $H_{kk} = \frac{d_k}{\hat{d}_k} \frac{F'_k \hat{F}_k}{T} = (1 + O_p(n^{-\frac{1}{2}\alpha_k}))(1 + \bar{O}_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k})) = 1 + O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k})$ by Lemma 7. Next, consider the case $\alpha_k \geq \alpha_l$. By Lemma 7

$$H_{lk} = \frac{d_l}{\hat{d}_k} \frac{F'_l \hat{F}_k}{T} = O_p(n^{\alpha_l-\alpha_k}) \bar{O}_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) = \bar{O}_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}). \quad (146)$$

Finally, from Lemma 8:

$$\frac{(\hat{F}'_k - F H_{\cdot k})' (\hat{F}'_k - F H_{\cdot k})}{T} = \frac{\hat{F}'_k \hat{F}_k}{T} + H'_{\cdot k} \frac{F' F}{T} H_{\cdot k} - 2 \frac{\hat{F}'_k F}{T} H_{\cdot k} = O_p(n^{1-2\alpha_k}). \quad (147)$$

Further,

$$\frac{\hat{F}'_k \hat{F}_k}{T} + H'_{\cdot k} \frac{F' F}{T} H_{\cdot k} - 2 \frac{\hat{F}'_k F}{T} H_{\cdot k} = 1 + \sum_{l=1}^r H_{lk}^2 - 2 \sum_{l=1}^r \frac{\hat{F}'_k F_l}{T} H_{lk} \quad (148)$$

$$= 1 + H_{kk}^2 - 2 \frac{\hat{F}'_k F_k}{T} H_{kk} + \sum_{l \neq k}^r H_{lk}^2 - 2 \sum_{l \neq k}^r \frac{\hat{F}'_k F_l}{T} H_{lk} \quad (149)$$

Since $H_{kk} = 1 + O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k})$, it follows that $1 + H_{kk}^2 - 2 \frac{\hat{F}'_k F_k}{T} H_{kk} = O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k})$. Combining this with (147)-(149), we obtain

$$\sum_{l \neq k}^r H_{lk}^2 - 2 \sum_{l \neq k}^r \frac{\hat{F}'_k F_l}{T} H_{lk} + O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) = O_p(n^{1-2\alpha_k}) \quad (150)$$

$$\sum_{l \neq k}^r \left(H_{lk}^2 - 2 \frac{\hat{F}'_k F_l}{T} H_{lk} \right) = O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) + O_p(n^{1-2\alpha_k}) \quad (151)$$

$$\sum_{l \neq k}^r \left(\left(\frac{d_l}{\hat{d}_k} \right)^2 \frac{F'_l \hat{F}_k}{T}^2 - 2 \left(\frac{\hat{F}'_k F_l}{T} \right)^2 \frac{d_l}{\hat{d}_k} \right) = O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) + O_p(n^{1-2\alpha_k}) \quad (152)$$

$$\sum_{l \neq k}^r \frac{d_l}{\hat{d}_k} \left(\frac{F'_l \hat{F}_k}{T} \right)^2 \left[\frac{d_l}{\hat{d}_k} - 2 \right] = O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) + O_p(n^{1-2\alpha_k}) \quad (153)$$

Split the sum above into three parts according to the relationship between α_k and α_l and start with the elements for which $\alpha_k > \alpha_l$. Then

$$\frac{d_l}{\hat{d}_k} \left(\frac{F'_l \hat{F}_k}{T} \right)^2 \left[\frac{d_l}{\hat{d}_k} - 2 \right] = O_p(n^{\alpha_l-\alpha_k}) O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) O_p(1) = o_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) \quad (154)$$

Next, consider elements in the sum for which $\alpha_k = \alpha_l$. Then

$$\frac{d_l}{\hat{d}_k} \left(\frac{F'_l \hat{F}_k}{T} \right)^2 \left[\frac{d_l}{\hat{d}_k} - 2 \right] = O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) O_p(1) = O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) \quad (155)$$

Finally consider the remaining terms. First note that for the remaining sum the upper limit for the entire sum still holds, as the terms in the first two cases are small enough. Further note that all terms in this remaining sum are positive with probability 1. Thus each term is bounded by their overall sum and for all k such that $\alpha_k < \alpha_l$:

$$\frac{d_l}{\hat{d}_k} \left(\frac{F'_l \hat{F}_k}{T} \right)^2 \left[\frac{d_l}{\hat{d}_k} - 2 \right] = O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) + O_p(n^{1-2\alpha_k}) \quad (156)$$

Since the LHS in (156) is equal to H_{lk}^2 up to a negligible term, this establishes that $H_{.k} = \iota_k + O_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2}-\alpha_k})$ in this last case, thus finishing the proof. \square

Lemma 10. *Under Assumptions 1-4, with \hat{F} and H defined as in the previous lemma:*

$$\hat{F}_{tk} - H'_{.k} F_t = O_p(n^{1-2\alpha_k}) + O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) \quad (157)$$

Proof. Revisit the decomposition from Lemma 8. It follows that:

$$\hat{F}_{tk} - H'_{.k} F_t = \hat{d}_k^{-1} \left(\frac{1}{T} \sum_{s=1}^T \hat{F}_{sk} e'_s e_t + \frac{1}{T} \sum_{s=1}^T \hat{F}_{sk} F'_s \Lambda' e_t + \frac{1}{T} \sum_{s=1}^T \hat{F}_{sk} e'_s \Lambda F_t \right) \quad (158)$$

$$= \hat{d}_k^{-1} \left(I_{kt} + II_{kt} + III_{kt}, \right) \quad (159)$$

Start with I_{kt} and decompose as follows:

$$I_{kt} = \frac{1}{T} \sum_{s=1}^T \hat{F}_{sk} e'_s e_t \quad (160)$$

$$\leq \frac{1}{T} \sum_{s=1}^T (\hat{F}_{sk} - H'_{.k} F_s) e'_s e_t + \frac{1}{T} H'_{.k} \sum_{s=1}^T F_s e'_s e_t \quad (161)$$

$$\begin{aligned} &\leq \frac{1}{T} \sum_{s=1}^T (\hat{F}_{sk} - H'_{.k} F_s) [e'_s e_t - \mathbb{E}(e'_s e_t)] + \frac{1}{T} H'_{.k} \sum_{s=1}^T F_s [e'_s e_t - \mathbb{E}(e'_s e_t)] \\ &\quad + \frac{1}{T} \sum_{s=1}^T (\hat{F}_{sk} - H'_{.k} F_s) \mathbb{E}(e'_s e_t) + \frac{1}{T} H'_{.k} \sum_{s=1}^T F_s \mathbb{E}(e'_s e_t) \end{aligned} \quad (162)$$

For the first part:

$$\begin{aligned} &\left\| \frac{1}{T} \sum_{s=1}^T (\hat{F}_{sk} - H'_{.k} F_s) [e'_s e_t - \mathbb{E}(e'_s e_t)] \right\| \\ &\leq \left(\frac{1}{T} \sum_{s=1}^T \|\hat{F}_{sk} - H'_{.k} F_s\|^2 \right)^{1/2} \left(\frac{1}{T} \sum_{s=1}^T [e'_s e_t - \mathbb{E}(e'_s e_t)]^2 \right)^{1/2} \end{aligned} \quad (163)$$

By Lemma 8 the first term is $\bar{O}_p(n^{\frac{1}{2}-\alpha_k})$. For the second term inside the brackets of (163):

$$\frac{1}{T} \sum_{s=1}^T [e'_s e_t - \mathbb{E}(e'_s e_t)]^2 = \frac{n}{T} \sum_{s=1}^T \left[\frac{1}{\sqrt{n}} e'_s e_t - \mathbb{E}(e'_s e_t) \right]^2. \quad (164)$$

This is $O_p(n)$ by Assumption 3(c), and thus the first part of the decomposition of I_t is $O_p(n^{\frac{1}{2}-\alpha_k}) O_p(\sqrt{n}) =$

$O_p(n^{1-\alpha_k})$. For the second part in the decomposition of I_{kt} :

$$H'_{k\cdot} \frac{1}{T} \sum_{s=1}^T F_s [e'_s e_t - \mathbb{E}(e'_s e_t)] = [\iota_k + O_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2}-\alpha_k})] O_p(1) \quad (165)$$

$$= O_p(1) + O_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2}-\alpha_k}), \quad (166)$$

by Assumption 3(d). Next consider the third part of I_{kt} :

$$\left| \frac{1}{T} \sum_{s=1}^T (\hat{F}_{sk} - H'_{k\cdot} F_s) \mathbb{E}(e'_s e_t) \right| \leq \left(\frac{1}{T} \sum_{s=1}^T (\hat{F}_{sk} - H'_{k\cdot} F_s)^2 \right)^{\frac{1}{2}} \frac{n}{\sqrt{T}} \left(\sum_{s=1}^T \mathbb{E}(\frac{e'_s e_t}{n})^2 \right)^{\frac{1}{2}} \quad (167)$$

$$= \frac{n}{\sqrt{T}} O_p(n^{\frac{1}{2}-\alpha_k}) O_p(1) = O_p(n^{1-\alpha_k}), \quad (168)$$

by Lemma 8 and Assumption 3(b). Finally, for the last part of I_{kt} , $\frac{1}{T} H'_{k\cdot} \sum_{s=1}^T F_s \mathbb{E}(e'_s e_t) = O_p(1) + O_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2}-\alpha_k})$, since

$$\mathbb{E} \left| \sum_{s=1}^T F_s E(e'_s e_t) \right| \leq \max_s \|F_s\| \sum_{s=1}^T |\mathbb{E}(e'_s e_t)| \leq C \quad (169)$$

by Assumption 3(b) and using the fact that $\max_s \|F_s\| < C$. It follows that

$$I_{kt} = O_p(n^{1-\alpha_k}) \quad (170)$$

Next, consider II_{kt} :

$$\begin{aligned} II_{kt} &= \frac{1}{T} \sum_{s=1}^T \hat{F}_{sk} F'_s \Lambda' e_t \\ &= \frac{1}{T} \sum_{s=1}^T (\hat{F}_{sk} - H'_{k\cdot} F_s) F'_s \Lambda' e_t + H'_{k\cdot} \frac{1}{T} \sum_{s=1}^T F_s F'_s \Lambda' e_t \end{aligned} \quad (171)$$

For the second part:

$$H'_{k\cdot} \frac{1}{T} \sum_{s=1}^T F_s F'_s \Lambda' e_t = H'_{k\cdot} \left(\frac{1}{T} \sum_{s=1}^T F_s F'_s \right) (\Lambda' e_t) = H'_{k\cdot} \Lambda' e_t \quad (172)$$

$$= [\iota_k + O_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2}-\alpha_k})] O_p(n^{\frac{1}{2}\alpha_{max}}) \quad (173)$$

$$= O_p(n^{\frac{1}{2}\alpha_{max}}) + O_p(n^{\frac{3}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2}+\frac{1}{2}\alpha_{max}-\alpha_k}). \quad (174)$$

For the first part:

$$\left\| \frac{1}{T} \sum_{s=1}^T (\hat{F}_{sk} - H'_{k \cdot} F_s) F'_s \Lambda' e_t \right\| \leq \left(\frac{1}{T} \sum_{s=1}^T \|(\hat{F}_{sk} - H'_{k \cdot} F_s)\|^2 \right)^{1/2} \left(\frac{1}{T} \sum_{s=1}^T \|F'_s \Lambda' e_t\|^2 \right)^{1/2}. \quad (175)$$

Further:

$$\frac{1}{T} \sum_{s=1}^T \|F'_s \Lambda' e_t\|^2 \leq \|\Lambda' e_t\|^2 \frac{1}{T} \sum_{s=1}^T \|F_s\|^2 = O_p(n^{\alpha_{max}}), \quad (176)$$

and by Lemma 8 $\left(\frac{1}{T} \sum_{s=1}^T \|(\hat{F}_{sk} - H'_{k \cdot} F_s)\|^2 \right)^{1/2} = O_p(n^{\frac{1}{2} - \alpha_k})$. Therefore:

$$\begin{aligned} II_{kt} &= \frac{1}{T} \sum_{s=1}^T \hat{F}_{sk} F'_s \Lambda' e_t \\ &= \frac{1}{T} \sum_{s=1}^T (\hat{F}_{sk} - H'_{k \cdot} F_s) F'_s \Lambda' e_t + H'_{k \cdot} \frac{1}{T} \sum_{s=1}^T F_s F'_s \Lambda' e_t \end{aligned} \quad (177)$$

$$\leq O_p(n^{\frac{1}{2} \alpha_{max}}) + O_p(n^{\frac{3}{4} \alpha_{max} - \frac{1}{2} \alpha_k}) + O_p(n^{\frac{1}{2} + \frac{1}{2} \alpha_{max} - \alpha_k}). \quad (178)$$

Finally, consider III_{kt} :

$$III_{kt} = \frac{1}{T} \sum_{s=1}^T \hat{F}_{sk} e'_s \Lambda F_t \quad (179)$$

$$= \frac{1}{T} \sum_{s=1}^T (\hat{F}_{sk} - H'_{k \cdot} F_s) e'_s \Lambda F_t + H'_{k \cdot} \frac{1}{T} \sum_{s=1}^T F_s e'_s \Lambda F_t. \quad (180)$$

Start with the first term:

$$\left\| \frac{1}{T} \sum_{s=1}^T (\hat{F}_{sk} - H'_{k \cdot} F_s) e'_s \Lambda F_t \right\| \leq \left(\frac{1}{T} \sum_{s=1}^T \|(\hat{F}_{sk} - H'_{k \cdot} F_s)\|^2 \right)^{1/2} \left(\frac{1}{T} \sum_{s=1}^T \|e'_s \Lambda\|^2 \right)^{1/2} \|F_t\| \quad (181)$$

$$= O_p(n^{\frac{1}{2} - \alpha_k}) O_p(n^{\frac{1}{2} \alpha_{max}}) O_p(1) = O_p(n^{\frac{1}{2} + \frac{1}{2} \alpha_{max} - \alpha_k}). \quad (182)$$

For the second term:

$$H'_{k \cdot} \frac{1}{T} \sum_{s=1}^T F_s e'_s \Lambda F_t = H'_{k \cdot} \frac{n^{\frac{1}{2} \alpha_{max}}}{\sqrt{T}} \left(\frac{1}{n^{\frac{1}{2} \alpha_{max}} \sqrt{T}} \sum_{s=1}^T F_s e'_s \Lambda \right) F_t \quad (183)$$

$$= [\iota_k + O_p(n^{\frac{1}{4} \alpha_{max} - \frac{1}{2} \alpha_k}) + O_p(n^{\frac{1}{2} - \alpha_k})] O_p(n^{\frac{1}{2} \alpha_{max} - \frac{1}{2}}) O_p(1) \quad (184)$$

$$= O_p(n^{\frac{1}{2} \alpha_{max} - \frac{1}{2}}) + O_p(n^{\frac{1}{2} \alpha_{max} - \alpha_k}) + O_p(n^{\frac{3}{4} \alpha_{max} - \frac{1}{2} - \frac{1}{2} \alpha_k}), \quad (185)$$

using Assumption 4(b). It follows that

$$III_{kt} = O_p(n^{\frac{1}{2} + \frac{1}{2}\alpha_{max} - \alpha_k}) + O_p(n^{\frac{1}{2}\alpha_{max} - \frac{1}{2}}) + O_p(n^{\frac{1}{2}\alpha_{max} - \alpha_k}) + O_p(n^{\frac{3}{4}\alpha_{max} - \frac{1}{2} - \frac{1}{2}\alpha_k}) \quad (186)$$

$$= O_p(n^{\frac{1}{2} + \frac{1}{2}\alpha_{max} - \alpha_k}) \quad (187)$$

Combining these partial results I obtain that:

$$\hat{F}_{tk} - H'_{k\cdot} F_t = \hat{d}_k^{-1} (I_t + II_t + III_t) \quad (188)$$

$$= O_p(n^{-\alpha_k}) \left(O_p(n^{1-\alpha_k}) + O_p(n^{\frac{1}{2}\alpha_{max}}) + O_p(n^{\frac{3}{4}\alpha_{max} - \frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2} + \frac{1}{2}\alpha_{max} - \alpha_k}) \right) \quad (189)$$

$$= O_p(n^{-\alpha_k}) \left(O_p(n^{1-\alpha_k}) + O_p(n^{\frac{1}{2}\alpha_{max}}) \right) \quad (190)$$

$$= O_p(n^{1-2\alpha_k}) + O_p(n^{\frac{1}{2}\alpha_{max} - \alpha_k}) \quad (191)$$

Note that we achieve convergence as long as $\alpha_k > .5$. Further, in the case of r strong factors ($\alpha_{min} = \alpha_{max} = 1$), (188) reduces to:

$$\hat{F}_t - H' F_t = O_p(n^{-1}) \left(O_p(1) + O_p(n^{\frac{1}{2}}) + O_p(1) \right) = O_p\left(\frac{1}{\sqrt{n}}\right) \quad (192)$$

This is in line with the literature (Bai (2003)). \square

Lemma 11. *Under Assumptions 1-4, with \hat{F} and H defined as in the previous lemmata:*

$$\frac{(\hat{F}_k - FH_{\cdot k})' F}{T} = \min\{O_p(n^{1-2\alpha_k}), O_p(n^{\frac{1}{2}-\alpha_k})\} \quad (193)$$

Proof. Using Identity (128),

$$\frac{(\hat{F}_k - FH_{\cdot k})' F}{T} = \frac{1}{T} \sum_{t=1}^T (\hat{F}_{tk} - H'_{k\cdot} F_t) F'_t \quad (194)$$

$$= \hat{d}_k^{-1} \left(\frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T \hat{F}_{sk} F'_t e'_s e_t + \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T \hat{F}_{sk} F'_t F'_s \Lambda' e_t + \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T \hat{F}_{sk} F'_t e'_s \Lambda F_t \right) \quad (195)$$

$$= \hat{d}_k^{-1} \left(I_k + II_k + III_k \right) \quad (196)$$

For I_k :

$$I_k = \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T \hat{F}_{sk} F'_t e'_s e_t \quad (197)$$

$$= \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T (\hat{F}_{sk} - H'_{k \cdot} F_s) F'_t e'_s e_t + \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T H'_{k \cdot} F_s F'_t e'_s e_t \quad (198)$$

$$\begin{aligned} &= \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T (\hat{F}_{sk} - H'_{k \cdot} F_s) F'_t [e'_s e_t - \mathbb{E}(e'_s e_t)] + \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T H'_{k \cdot} F_s F'_t [e'_s e_t - \mathbb{E}(e'_s e_t)] \\ &\quad + \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T (\hat{F}_{sk} - H'_{k \cdot} F_s) F'_t \mathbb{E}(e'_s e_t) + \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T H'_{k \cdot} F_s F'_t \mathbb{E}(e'_s e_t) \end{aligned} \quad (199)$$

The same arguments as in the proof of Lemma B.2 in Bai (2003) can be used to show that $I_k = O_p(n^{1-\alpha_k})$. Details are omitted. For II_k :

$$II_k = \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T \hat{F}_{sk} F'_t F'_s \Lambda' e_t \quad (200)$$

$$= \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T (\hat{F}_{sk} - H'_{k \cdot} F_s) F'_t F'_s \Lambda' e_t + \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T H'_{k \cdot} F_s F'_t F'_s \Lambda' e_t \quad (201)$$

Consider both parts in turn:

$$\left\| \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T (\hat{F}_{sk} - H'_{k \cdot} F_s) F'_t F'_s \Lambda' e_t \right\| \quad (202)$$

$$= \left(\frac{1}{T} \sum_{s=1}^T \|\hat{F}_{sk} - H'_{k \cdot} F_s\|^2 \right)^{\frac{1}{2}} \left(\frac{1}{T} \sum_{s=1}^T \left\| \frac{1}{T} \sum_{t=1}^T F'_t F'_s \Lambda' e_t \right\|^2 \right)^{\frac{1}{2}} \quad (203)$$

$$= \left(\frac{1}{T} \sum_{s=1}^T \|\hat{F}_{sk} - H'_{k \cdot} F_s\|^2 \right)^{\frac{1}{2}} \left(\frac{1}{T} \sum_{s=1}^T (\|F_s\|^2 \left\| \frac{1}{T} \sum_{t=1}^T F'_t \Lambda' e_t \right\|^2) \right)^{\frac{1}{2}} \quad (204)$$

$$\leq O_p(n^{\frac{1}{2}-\alpha_k}) \frac{n^{\frac{1}{2}\alpha_{max}}}{\sqrt{T}} \left(\frac{1}{T} \sum_{s=1}^T (\|F_s\|^2 \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T F'_t \frac{\Lambda' e_t}{n^{\frac{1}{2}\alpha_{max}}} \right\|^2) \right)^{\frac{1}{2}} \quad (205)$$

$$= O_p(n^{\frac{1}{2}-\alpha_k}) O_p(n^{\frac{1}{2}\alpha_{max}-\frac{1}{2}}) O_p(1) \quad (206)$$

$$= O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}), \quad (207)$$

where the boundedness of the third term in (206) follows from Assumption 4(a). Further:

$$\frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T H'_{k \cdot} F_s F'_t F'_s \Lambda' e_t = H'_{k \cdot} \left(\frac{1}{T} \sum_{s=1}^T F_s F'_s \right) \frac{1}{T} \sum_{t=1}^T F'_t \Lambda' e_t \quad (208)$$

$$= [\iota_k + O_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2}-\alpha_k})] O_p(1) O_p(n^{\frac{1}{2}\alpha_{max}-\frac{1}{2}}) \quad (209)$$

$$= O_p(n^{\frac{1}{2}\alpha_{max}-\frac{1}{2}}) + O_p(n^{\frac{3}{4}\alpha_{max}-\frac{1}{2}-\frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) \quad (210)$$

$$= O_p(n^{\frac{1}{2}\alpha_{max}-\frac{1}{2}}) + O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}), \quad (211)$$

where the last rate again follows from Assumption 4(b). I conclude:

$$II = O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) + O_p(n^{\frac{1}{2}\alpha_{max}-\frac{1}{2}}) + O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) \quad (212)$$

$$= O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) + O_p(n^{\frac{1}{2}\alpha_{max}-\frac{1}{2}}) \quad (213)$$

The proof of III follows the same arguments as that of II. Combining rates, I therefore conclude that

$$\frac{(\hat{F}_k - FH_{\cdot k})' F}{T} = \hat{d}_k^{-1} \left(I_k + II_k + III_k \right) \quad (214)$$

$$= O_p(n^{-\alpha_k}) \left(O_p(n^{1-\alpha_k}) + O_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) + O_p(n^{\frac{1}{2}\alpha_{max}-\frac{1}{2}}) \right) \quad (215)$$

$$= O_p(n^{1-2\alpha_k}), \quad (216)$$

as the latter terms are dominated by the former. \square

Lemma 12. Under Assumptions 1-4, with \hat{F} and H defined as in the previous lemmata:

$$\frac{(\hat{F}_k - FH_{\cdot k})' e_i}{T} = O_p(n^{1-2\alpha_k}) \quad (217)$$

Proof.

$$\frac{(\hat{F}_k - FH_{\cdot k})' e_i}{T} = \frac{1}{T} \sum_{t=1}^T (\hat{F}_{tk} - H'_{k \cdot} F_t) e_{it} \quad (218)$$

$$= \hat{d}_k^{-1} \left(\frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T \hat{F}_{sk} e'_s e_t e_{it} + \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T \hat{F}_{sk} F'_s \Lambda' e_t e_{it} + \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T \hat{F}_{sk} e'_s \Lambda F_t e_{it} \right) \quad (219)$$

$$= \hat{d}_k^{-1} \left(I_k + II_k + III_k \right) \quad (220)$$

$$I_k = \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T \hat{F}_{sk} e'_s e_t e_{it} \quad (221)$$

$$= \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T (\hat{F}_{sk} - H'_{k \cdot} F_s) e'_s e_t e_{it} + \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T H'_{k \cdot} F_s e'_s e_t e_{it} \quad (222)$$

$$\begin{aligned} &= \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T (\hat{F}_{sk} - H'_{k \cdot} F_s) [e'_s e_t - \mathbb{E}(e'_s e_t)] e_{it} + \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T H'_{k \cdot} F_s [e'_s e_t - \mathbb{E}(e'_s e_t)] e_{it} \\ &\quad + \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T (\hat{F}_{sk} - H'_{k \cdot} F_s) \mathbb{E}(e'_s e_t) e_{it} + \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T H'_{k \cdot} F_s \mathbb{E}(e'_s e_t) e_{it} \end{aligned} \quad (223)$$

Consider these four terms in turn.

$$\begin{aligned} &\frac{1}{T^2} \sum_{s=1}^T \sum_{t=1}^T (\hat{F}_{sk} - H'_{k \cdot} F_s) [e'_s e_t - \mathbb{E}(e'_s e_t)] e_{it} \\ &\leq \left(\frac{1}{T} \sum_{s=1}^T \|\hat{F}_{sk} - H'_{k \cdot} F_s\|^2 \right)^{\frac{1}{2}} \left(\frac{1}{T} \sum_{s=1}^T \left(\frac{1}{T} \sum_{t=1}^T [e'_s e_t - \mathbb{E}(e'_s e_t)] e_{it} \right)^2 \right)^{\frac{1}{2}} \end{aligned} \quad (224)$$

$$\leq \sqrt{n} \left(\frac{1}{T} \sum_{s=1}^T \|\hat{F}_{sk} - H'_{k \cdot} F_s\|^2 \right)^{\frac{1}{2}} \left(\frac{1}{T} \sum_{s=1}^T \left(\frac{1}{T} \sum_{t=1}^T \frac{1}{\sqrt{n}} [e'_s e_t - \mathbb{E}(e'_s e_t)] e_{it} \right)^2 \right)^{\frac{1}{2}} \quad (225)$$

$$\leq \sqrt{n} O_p(n^{\frac{1}{2}-\alpha_k}) O_p(1) = O_p(n^{1-\alpha_k}), \quad (226)$$

where the boundedness of the last term follows from Assumption 3(c). For the next term, ignoring H , take expectations:

$$\mathbb{E} \left[\frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T F_s [e'_s e_t - \mathbb{E}(e'_s e_t)] e_{it} \right] = \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \left(\frac{1}{\sqrt{nT}} \sum_{s=1}^T F_s [e'_s e_t - \mathbb{E}(e'_s e_t)] \right) e_{it} \right] \quad (227)$$

$$\leq \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left(\left\| \frac{1}{\sqrt{nT}} \sum_{s=1}^T F_s [e'_s e_t - \mathbb{E}(e'_s e_t)] \right\|^2 \right)^{\frac{1}{2}} (\mathbb{E}(e_{it})^2)^{\frac{1}{2}} \quad (228)$$

$$= O(1) \quad (229)$$

For the third term:

$$\begin{aligned} &\frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T (\hat{F}_{sk} - H'_{k \cdot} F_s) \mathbb{E}(e'_s e_t) e_{it} \\ &\leq \left(\frac{1}{T} \sum_{s=1}^T \|\hat{F}_{sk} - H'_{k \cdot} F_s\|^2 \right)^{\frac{1}{2}} \left(\frac{n}{T} \sum_{t=1}^T \sum_{s=1}^T |\mathbb{E}(\frac{e'_s e_t}{n})|^2 \frac{1}{T} \sum_{t=1}^T e_{it}^2 \right)^{\frac{1}{2}} \end{aligned} \quad (230)$$

$$= O_p(n^{\frac{1}{2}-\alpha_k}) O_p(\sqrt{n}) = O_p(n^{1-\alpha_k}), \quad (231)$$

using Lemma 3. Finally, ignoring H, take expectations of the last term:

$$\mathbb{E} \left[\frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T F_s \mathbb{E}(e'_s e_t) e_{it} \right] \leq \frac{1}{T} \sum_{t=1}^T \sum_{s=1}^T (\|F_s\|^2)^{\frac{1}{2}} \mathbb{E}\left(\frac{e'_s e_t}{n}\right) (\mathbb{E} e_{it}^2)^{\frac{1}{2}} = O(1), \quad (232)$$

since both the first and third term in the final sum are bounded and using Assumption 3(b). Therefore $I_k = O_p(n^{1-\alpha_k}) + O_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2}-\alpha_k}) = O_p(n^{1-\alpha_k})$. Next consider II_k :

$$II_k = \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T \hat{F}_{sk} F'_s \Lambda' e_t e_{it} \quad (233)$$

$$= \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T (\hat{F}_{sk} - H'_{k \cdot} F_s) F'_s \Lambda' e_t e_{it} + \frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T H'_{k \cdot} F_s F'_s \Lambda' e_t e_{it}. \quad (234)$$

Again consider both terms separately and start with the second:

$$\frac{1}{T^2} \sum_{t=1}^T \sum_{s=1}^T H'_{k \cdot} F_s F'_s \Lambda' e_t e_{it} = H'_{k \cdot} \left(\frac{1}{T} \sum_{s=1}^T F_s F'_s \right) \frac{1}{T} \sum_{t=1}^T \sum_{j=1}^n \lambda_j e_{jt} e_{it} \quad (235)$$

$$= H'_{k \cdot} \left(\frac{1}{T} \sum_{s=1}^T F_s F'_s \right) \left(\frac{1}{T} \sum_{t=1}^T \sum_{j=1}^n \lambda_j [e_{jt} e_{it} - \mathbb{E}(e_{jt} e_{it})] + \frac{1}{T} \sum_{t=1}^T \sum_{j=1}^n \lambda_j \mathbb{E}(e_{jt} e_{it}) \right) \quad (236)$$

$$\leq H'_{k \cdot} \left(\frac{1}{T} \sum_{s=1}^T F_s F'_s \right) \left(\frac{C}{\sqrt{T} \sqrt{n}} \sum_{t=1}^T \sum_{j=1}^n [e_{jt} e_{it} - \mathbb{E}(e_{jt} e_{it})] + \sum_{j=1}^n \lambda_j \mathbb{E}\left(\frac{e'_j e_i}{T}\right) \right) \quad (237)$$

$$= [\iota_k + O_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2}-\alpha_k})][O_p(1) + O(1)], \quad (238)$$

where the boundedness of the last term follows from Assumption 3(b). Similarly for the first term:

$$\frac{1}{T^2} \sum_{s=1}^T \sum_{t=1}^T (\hat{F}_{sk} - H'_{k \cdot} F_s) F'_s \Lambda' e_t e_{it} \quad (239)$$

$$\leq \left(\frac{1}{T} \sum_{s=1}^T \|\hat{F}_{sk} - H'_{k \cdot} F_s\|^2 \right)^{\frac{1}{2}} \left(\frac{1}{T} \sum_{s=1}^T \left(\frac{1}{T} \sum_{t=1}^T F'_s \Lambda' e_t e_{it} \right)^2 \right)^{\frac{1}{2}} \quad (240)$$

$$\leq \left(\frac{1}{T} \sum_{s=1}^T \|\hat{F}_{sk} - H'_{k \cdot} F_s\|^2 \right)^{\frac{1}{2}} \left(\frac{1}{T} \sum_{s=1}^T \left(F'_s \frac{1}{T} \sum_{t=1}^T \sum_{j=1}^n \lambda_j e_{jt} e_{it} \right)^2 \right)^{\frac{1}{2}} \quad (241)$$

$$= O_p(n^{\frac{1}{2}-\alpha_k}) O_p(1), \quad (242)$$

using the same arguments as above. We conclude that $II_k = O_p(1) + O_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2}-\alpha_k})$. Finally, using similar arguments as in the proof of II_k , one can show that the same bounds apply

to III_k , and it follows that

$$\frac{(\hat{F}_k - FH_{\cdot k})' e_i}{T} = \hat{d}_k^{-1} \left(I_k + II_k + III_k \right) \quad (243)$$

$$= O_p(n^{-\alpha_k}) \left(O_p(n^{1-\alpha_k}) + O_p(1) + O_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + O_p(n^{\frac{1}{2}-\alpha_k}) \right) \quad (244)$$

$$= O_p(n^{1-2\alpha_k}) \quad (245)$$

□

Lemma 13. Under Assumptions 1-4, let $z = n^\tau \sqrt{\log(n)}$ for some $\tau \in [.5, 1]$. With slight abuse of notation, the estimated loadings $\hat{\lambda}_{ik}$ are ordered such that, for each k , $|\hat{\lambda}_{1k}| \geq |\hat{\lambda}_{2k}| \geq \dots \geq |\hat{\lambda}_{nk}|$. Then

- (a) If $\alpha_k > \max\{\frac{1+\tau}{3}, \frac{\alpha_{max}+4\tau}{6}\}$: $\frac{1}{z} \sum_{i=1}^z \hat{\lambda}_{ik}^2 - \frac{1}{z} \sum_{i=1}^z \lambda_{ik}^2 = \bar{O}_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + \bar{O}_p(n^{1-2\alpha_k})$
- (b) If $\alpha_k \leq \max\{\frac{1+\tau}{3}, \frac{\alpha_{max}+4\tau}{6}\}$: $\frac{1}{z} \sum_{i=1}^z \hat{\lambda}_{ik}^2 - \frac{1}{z} \sum_{i=1}^z \lambda_{ik}^2 = O_p\left(\frac{n^{\alpha_k}}{n^\tau \sqrt{\log(n)}}\right)$

Proof. By Theorem 3:

$$\hat{\lambda}_{ik} - \lambda_{ik} = \bar{O}_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + \bar{O}_p(n^{1-2\alpha_k}). \quad (246)$$

Since

$$\frac{1}{z} \sum_{i=1}^z \hat{\lambda}_{ik}^2 - \frac{1}{z} \sum_{i=1}^z \lambda_{ik}^2 = \frac{1}{z} \sum_{i=1}^z (\hat{\lambda}_{ik}^2 - \lambda_{ik}^2), \quad (247)$$

this is just an average (squared) deviation and the result in part (a) immediately follows.

Next consider the case $\alpha_k \leq \tau$.

$$\frac{1}{z} \sum_{i=1}^z \lambda_{ik}^2 \leq \frac{1}{z} \sum_{i=1}^n \lambda_{ik}^2 = \frac{n^{-\tau}}{\sqrt{\log(n)}} \psi_k(\Lambda' \Lambda) = O_p\left(\frac{n^{\alpha_k-\tau}}{\sqrt{\log(n)}}\right) \quad (248)$$

and, similarly

$$\frac{1}{z} \sum_{i=1}^z \hat{\lambda}_{ik}^2 \leq \frac{1}{z} \sum_{i=1}^n \hat{\lambda}_{ik}^2 = \frac{n^{-\tau}}{\sqrt{\log(n)}} \psi_k\left(\frac{X' X}{T}\right) = O_p\left(\frac{n^{\alpha_k-\tau}}{\sqrt{\log(n)}}\right). \quad (249)$$

Combined they imply the stated bound on the difference.

Finally, consider the case $\max\{\frac{1+\tau}{3}, \frac{\alpha_{max}+4\tau}{6}\} < \alpha_k \leq \tau$. In those situations both bounds above apply and imply convergence to zero. For $\alpha_k > \max\{\frac{1+\tau}{3}, \frac{\alpha_{max}+4\tau}{6}\}$, the first bound is the tighter one and thus applies. \square

Lemma 14. *There exists a constant $c > 0$ such that $\lim_{n \rightarrow \infty} P(\hat{T}_{zk}^u / n^{(1+\frac{1}{2}u)\alpha_k - \frac{1}{2}u} < c) = 0$ for $k = 1, \dots, r_{max}$.*

Proof. First note that

$$\frac{1}{z} \sum_{ik}^z \hat{\lambda}_{ik}^2 \geq \frac{1}{n} \sum_{ik}^n \hat{\lambda}_{ik}^2 = \frac{1}{n} \psi_k \left(\frac{X'X}{T} \right) \quad (250)$$

It follows that

$$\hat{T}_{zk}^u = \psi_k \left(\frac{X'X}{T} \right) \left(\frac{1}{z} \sum_i^z \frac{\hat{\lambda}_{ik}^2}{\sqrt{\frac{1}{n} \sum_{i=1}^n \hat{\lambda}_{ik}^2}} \right)^u \geq \psi_k \left(\frac{X'X}{T} \right) \left(\frac{1}{z} \sum_i^z \hat{\lambda}_{ik}^2 \right)^{\frac{1}{2}u} \quad (251)$$

$$\geq \psi_k \left(\frac{X'X}{T} \right) \left[\frac{1}{n} \psi_k \left(\frac{X'X}{T} \right) \right]^{\frac{1}{2}u} = n^{-\frac{1}{2}u} \psi_k \left(\frac{X'X}{T} \right)^{1+\frac{1}{2}u} \quad (252)$$

For $\alpha_k > 0$, there exists a $c_1 > 0$ such that

$$\lim_{n \rightarrow \infty} P \left(\psi_k \left(\frac{X'X}{T} \right) / n^{\alpha_k} < c_1 \right) = 0 \quad (253)$$

and thus

$$\lim_{n \rightarrow \infty} P \left(\hat{T}_{zk}^u / n^{(1+\frac{1}{2}u)\alpha_k - \frac{1}{2}u} \geq c \right) = 1 \quad (254)$$

Finally, if $\alpha_k = 0$, since $\psi_k \left(\frac{X'X}{T} \right) > c_{eig} > 0$ for $k = r + 1, \dots, [dn]$ this implies that there exists a positive constant c_2 such that

$$\hat{T}_{zk}^u \geq c_2 n^{-\frac{1}{2}u} \quad (255)$$

\square

D.2 Proofs of Theorem 4 and Corollary 4

Proof of Corollary 4. First consider $k = r + 1, \dots, r_{max}$. Then, by Theorem 1, $\psi_k \left(\frac{XX'}{T} \right) = O_p(1)$ and thus there exists a finite $c_1 > 0$, $\lim_{n \rightarrow \infty} P \left(\hat{T}_k^2 \geq c_1 \right) = 0$. Further, by (66) following

Assumption 3 (e), there exists a constant $c_2 > 0$, such that $P\left(\psi_k\left(\frac{XX'}{T}\right) \geq c_2\right) = 1$ for $k^* = r+1, \dots, r_{max}$. Then, for any finite $c_3 > 0$,

$$\begin{aligned} & \lim_{n \rightarrow \infty} P\left(\frac{\psi_k\left(\frac{XX'}{T}\right)}{\psi_{k+1}\left(\frac{XX'}{T}\right)} > c_3 \log(n)\right) \\ &= \lim_{n \rightarrow \infty} \left[P\left(\frac{\psi_k\left(\frac{XX'}{T}\right)}{\psi_{k+1}\left(\frac{XX'}{T}\right)} > c_3 \log(n) \middle| \psi_{k+1}\left(\frac{XX'}{T}\right) < c_2\right) P\left(\psi_{k+1}\left(\frac{XX'}{T}\right) < c_2\right) \right. \\ &\quad \left. + P\left(\frac{\psi_k\left(\frac{XX'}{T}\right)}{\psi_{k+1}\left(\frac{XX'}{T}\right)} > c_3 \log(n) \middle| \psi_{k+1}\left(\frac{XX'}{T}\right) \geq c_2\right) P\left(\psi_{k+1}\left(\frac{XX'}{T}\right) \geq c_2\right) \right] \quad (256) \end{aligned}$$

$$= \lim_{n \rightarrow \infty} P\left(\frac{\psi_k\left(\frac{XX'}{T}\right)}{\psi_{k+1}\left(\frac{XX'}{T}\right)} > c_3 \log(n) \middle| \psi_{k+1}\left(\frac{XX'}{T}\right) \geq c_2\right) + 0 \quad (257)$$

$$\leq \lim_{n \rightarrow \infty} P\left(\psi_k\left(\frac{XX'}{T}\right) > c_2 c_3 \log(n)\right) = 0. \quad (258)$$

Next, consider $k = 1, \dots, r-1$. We already established for any finite $q_1 > 0$ that $\lim_{n \rightarrow \infty} P\left(\psi_k\left(\frac{XX'}{T}\right) > q_1\sqrt{n}\right)$

1. It then immediately follows that, there exists an $h > 0$ such that

$$\begin{aligned} & \lim_{n \rightarrow \infty} P\left(\frac{\psi_k\left(\frac{XX'}{T}\right)}{\psi_{k+1}\left(\frac{XX'}{T}\right)} > h\sqrt{n}\right) \\ &= \lim_{n \rightarrow \infty} \left[P\left(\frac{\psi_k\left(\frac{XX'}{T}\right)}{\psi_{k+1}\left(\frac{XX'}{T}\right)} > h\sqrt{n} \middle| \psi_{k+1}\left(\frac{XX'}{T}\right) < q_1\sqrt{n}\right) P\left(\psi_{k+1}\left(\frac{XX'}{T}\right) < q_1\sqrt{n}\right) \right. \\ &\quad \left. + P\left(\frac{\psi_r\left(\frac{XX'}{T}\right)}{\psi_{k+1}\left(\frac{XX'}{T}\right)} > h\sqrt{n} \middle| \psi_{k+1}\left(\frac{XX'}{T}\right) \geq q_1\sqrt{n}\right) P\left(\psi_{r+1}\left(\frac{XX'}{T}\right) \geq q_1\sqrt{n}\right) \right] \quad (259) \end{aligned}$$

$$= \lim_{n \rightarrow \infty} P\left(\frac{\psi_k\left(\frac{XX'}{T}\right)}{\psi_{k+1}\left(\frac{XX'}{T}\right)} \geq h\sqrt{n} \middle| \psi_{k+1}\left(\frac{XX'}{T}\right) \geq q_1\sqrt{n}\right) + 0 \quad (260)$$

$$\leq \lim_{n \rightarrow \infty} P\left(\psi_r\left(\frac{XX'}{T}\right) > q_1 hn\right). \quad (261)$$

But since there exists a finite $q_2 > 0$ with $\lim_{n \rightarrow \infty} P\left(\psi_k\left(\frac{XX'}{T}\right) > q_2 n\right) = 0$, letting $h = \frac{q_2}{q_1}$

establishes $\lim_{n \rightarrow \infty} P\left(\frac{\psi_k\left(\frac{XX'}{T}\right)}{\psi_{k+1}\left(\frac{XX'}{T}\right)} > h\sqrt{n}\right) = 0$. Finally, consider $k = r$. By Assumption 5

$\alpha_k > .5$ and thus $\lim_{n \rightarrow \infty} P\left(\psi_k\left(\frac{XX'}{T}\right) > q_1\sqrt{n}\right) = 1$ for any finite $q_1 > 0$. On the other hand,

$\psi_{r+1}\left(\frac{XX'}{T}\right) = O_p(1)$ and thus there exists a $q_2 > 0$, such that $P\left(\psi_{r+1}\left(\frac{XX'}{T}\right) \geq q_2\right) = 0$. Then,

for any finite $q_3 > 0$

$$\begin{aligned}
& \lim_{n \rightarrow \infty} P \left(\frac{\psi_r \left(\frac{XX'}{T} \right)}{\psi_{r+1} \left(\frac{XX'}{T} \right)} > q_3 \sqrt{n} \right) \\
&= \lim_{n \rightarrow \infty} \left[P \left(\frac{\psi_r \left(\frac{XX'}{T} \right)}{\psi_{r+1} \left(\frac{XX'}{T} \right)} > q_3 \sqrt{n} \middle| \psi_{r+1} \left(\frac{XX'}{T} \right) < q_2 \right) P \left(\psi_{r+1} \left(\frac{XX'}{T} \right) < q_2 \right) \right. \\
&\quad \left. + P \left(\frac{\psi_r \left(\frac{XX'}{T} \right)}{\psi_{r+1} \left(\frac{XX'}{T} \right)} > q_3 \sqrt{n} \middle| \psi_{r+1} \left(\frac{XX'}{T} \right) \geq q_2 \right) P \left(\psi_{r+1} \left(\frac{XX'}{T} \right) \geq q_2 \right) \right] \quad (262)
\end{aligned}$$

$$\lim_{n \rightarrow \infty} P \left(\frac{\psi_r \left(\frac{XX'}{T} \right)}{\psi_{r+1} \left(\frac{XX'}{T} \right)} > q_3 \sqrt{n} \middle| \psi_{r+1} \left(\frac{XX'}{T} \right) < q_2 \right) + 0 \quad (263)$$

$$\geq \lim_{n \rightarrow \infty} P \left(\psi_r \left(\frac{XX'}{T} \right) > q_2 q_3 \sqrt{n} \right) = 1 \quad (264)$$

Choosing $q_3 = \frac{h}{q_2}$, this completes the proof. \square

Proof of Theorem 4. First note that:

$$\hat{T}_{zk}^u - T_{zk}^u = n^{\frac{1}{2}u} \left[\psi_k \left(\frac{X'X}{T} \right)^{1-\frac{1}{2}u} \left(\frac{1}{z} \sum_i^z \hat{\lambda}_{ik}^2 \right)^u - \psi_k \left(\frac{\Lambda F' F \Lambda'}{T} \right)^{1-\frac{1}{2}u} \left(\frac{1}{z} \sum_i^z \lambda_{ik}^2 \right)^u \right] \quad (265)$$

Because $ab - cd = (a - c)d + (b - d)c + (a - c)(b - d)$ I may write

$$\hat{T}_{zk}^u - T_{zk}^u = n^{\frac{1}{2}u} \left[I + II + III \right], \quad (266)$$

where

$$I = \left(\psi_k \left(\frac{X'X}{T} \right)^{1-\frac{1}{2}u} - \psi_k \left(\frac{\Lambda F' F \Lambda'}{T} \right)^{1-\frac{1}{2}u} \right) \left(\frac{1}{z} \sum_i^z \lambda_{ik}^2 \right)^u \quad (267)$$

$$II = \left(\left(\frac{1}{z} \sum_i^z \hat{\lambda}_{ik}^2 \right)^u - \left(\frac{1}{z} \sum_i^z \lambda_{ik}^2 \right)^u \right) \psi_k \left(\frac{\Lambda F' F \Lambda'}{T} \right)^{1-\frac{1}{2}u} \quad (268)$$

$$III = \left(\psi_k \left(\frac{X'X}{T} \right)^{1-\frac{1}{2}u} - \psi_k \left(\frac{\Lambda F' F \Lambda'}{T} \right)^{1-\frac{1}{2}u} \right) \left(\left(\frac{1}{z} \sum_i^z \hat{\lambda}_{ik}^2 \right)^u - \left(\frac{1}{z} \sum_i^z \lambda_{ik}^2 \right)^u \right) \quad (269)$$

First consider the difference in I :

$$\psi_k \left(\frac{X'X}{T} \right)^{1-\frac{1}{2}u} - \psi_k \left(\frac{\Lambda F' F \Lambda'}{T} \right)^{1-\frac{1}{2}u} \quad (270)$$

$$= n^{(1-\frac{1}{2}u)\alpha_k} \left[\left(\psi_k \left(\frac{\Lambda F' F \Lambda'}{Tn^{\alpha_k}} \right) + \psi_k \left(\frac{X'X}{Tn^{\alpha_k}} \right) - \psi_k \left(\frac{\Lambda F' F \Lambda'}{Tn^{\alpha_k}} \right) \right)^{1-\frac{1}{2}u} - \psi_k \left(\frac{\Lambda F' F \Lambda'}{Tn^{\alpha_k}} \right)^{1-\frac{1}{2}u} \right] \quad (271)$$

$$= n^{(1-\frac{1}{2}u)\alpha_k} \left[\left(\psi_k \left(\frac{\Lambda F' F \Lambda'}{Tn^{\alpha_k}} \right) + \varepsilon_{\psi_k} \right)^{1-\frac{1}{2}u} - \psi_k \left(\frac{\Lambda F' F \Lambda'}{Tn^{\alpha_k}} \right)^{1-\frac{1}{2}u} \right], \quad (272)$$

where $\varepsilon_{\psi_k} = O_p(n^{-\frac{1}{2}\alpha_k})$ by Theorem 1. Using Newton's generalised binomial theorem:

$$\left(\psi_k \left(\frac{\Lambda F' F \Lambda'}{Tn^{\alpha_k}} \right) + \varepsilon_{\psi_k} \right)^{1-\frac{1}{2}u} - \psi_k \left(\frac{\Lambda F' F \Lambda'}{Tn^{\alpha_k}} \right)^{1-\frac{1}{2}u} \quad (273)$$

$$= \sum_{w=0}^{\infty} \frac{\Gamma(2-\frac{1}{2}u)}{w!} \psi_k \left(\frac{\Lambda F' F \Lambda'}{Tn^{\alpha_k}} \right)^{1-\frac{1}{2}u-w} \varepsilon_{\psi_k}^w - \psi_k \left(\frac{\Lambda F' F \Lambda'}{Tn^{\alpha_k}} \right)^{1-\frac{1}{2}u} \quad (274)$$

$$= O_p(n^{-\frac{1}{2}\alpha_k}) + o_p(n^{-\frac{1}{2}\alpha_k}) \quad (275)$$

I can thus distinguish between two cases as follows:

For $\alpha_k > \tau$: $I = n^{(1-\frac{1}{2}u)\alpha_k} [O_p(n^{-\frac{1}{2}\alpha_k}) + o_p(n^{-\frac{1}{2}\alpha_k})] O_p(1) = O_p(n^{(\frac{1}{2}-\frac{1}{2}u)\alpha_k})$

For $\alpha_k \leq \tau$: $I = n^{(1-\frac{1}{2}u)\alpha_k} [O_p(n^{-\frac{1}{2}\alpha_k}) + o_p(n^{-\frac{1}{2}\alpha_k})] O_p \left(\frac{n^{\alpha_k}}{n^\tau \log(n)^{\frac{1}{2}u}} \right) = O_p(n^{(\frac{1}{2}+\frac{1}{2}u)\alpha_k - \tau u} \log(n)^{-\frac{1}{2}u})$

Next, consider the difference in II . For (a) ($\alpha_k > \tau$):

$$\left(\frac{1}{z} \sum_i^z \hat{\lambda}_{ik}^2\right)^u - \left(\frac{1}{z} \sum_i^z \lambda_{ik}^2\right)^u \quad (276)$$

$$= \left(\frac{1}{z} \sum_i^z \lambda_{ik}^2 + \frac{1}{z} \sum_i^z \hat{\lambda}_{ik}^2 - \frac{1}{z} \sum_i^z \lambda_{ik}^2\right)^u - \left(\frac{1}{z} \sum_i^z \lambda_{ik}^2\right)^u \quad (277)$$

$$= \left(\frac{1}{z} \sum_i^z \lambda_{ik}^2 + \frac{1}{z} \sum_i^z (\hat{\lambda}_{ik}^2 - \lambda_{ik}^2)\right)^u - \left(\frac{1}{z} \sum_i^z \lambda_{ik}^2\right)^u \quad (278)$$

$$\begin{aligned} &= u \left[\frac{1}{z} \sum_i^z \lambda_{ik}^2 \right]^{u-1} \left(\frac{1}{z} \sum_i^z (\hat{\lambda}_{ik}^2 - \lambda_{ik}^2) \right) \\ &\quad + \mathbf{1}_{\{u>1\}} \frac{u(u-1)}{2} \left[\left(\frac{1}{z} \sum_i^z \lambda_{ik}^2 \right)^{u-2} \left(\frac{1}{z} \sum_i^z (\hat{\lambda}_{ik}^2 - \lambda_{ik}^2) \right)^2 \right] + \dots, \end{aligned} \quad (279)$$

where the third equality follows from the generalised binomial theorem for nonnegative exponents. Later terms will be dominated.

$$\begin{aligned} II &= [\mathbf{1}_{\{u>0\}} [\bar{O}_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + \bar{O}_p(n^{1-2\alpha_k})] \\ &\quad + \mathbf{1}_{\{u>1\}} [\bar{O}_p(n^{\frac{1}{2}\alpha_{max}-\alpha_k}) + \bar{O}_p(n^{2-4\alpha_k})]] O_p(n^{(1-\frac{1}{2}u)\alpha_k}) \end{aligned} \quad (280)$$

$$= \mathbf{1}_{\{u>0\}} [\bar{O}_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) + \bar{O}_p(n^{1-2\alpha_k})] O_p(n^{(1-\frac{1}{2}u)\alpha_k}) \quad (281)$$

Similarly, by Lemma 13, the same rate holds if $\max\{\frac{1+\tau}{3}, \frac{\alpha_{max}+4\tau}{6}\} < \alpha_k \leq \tau$. On the other hand, if $\alpha_k \leq \max\{\frac{1+\tau}{3}, \frac{\alpha_{max}+4\tau}{6}\}$, by Lemma 13:

$$\left(\frac{1}{z} \sum_i^z \hat{\lambda}_{ik}^2\right)^u - \left(\frac{1}{z} \sum_i^z \lambda_{ik}^2\right)^u = \frac{1}{z^u} \left[\left(\sum_i^z \hat{\lambda}_{ik}^2 \right)^u - \left(\sum_i^z \lambda_{ik}^2 \right)^u \right] \quad (282)$$

$$= \frac{1}{n^{\tau u} \log(n)^{\frac{1}{2}u}} \mathbf{1}_{\{u>0\}} [O_p(n^{u\alpha_k}) - O_p(n^{u\alpha_k})] \quad (283)$$

$$= \mathbf{1}_{\{u>0\}} O_p(n^{(\alpha_k-\tau)u} \log(n)^{-\frac{1}{2}u}), \quad (284)$$

which in turn implies that

$$II = \mathbf{1}_{\{u>0\}} O_p\left(\frac{n^{(\alpha_k-\tau)u}}{\log(n)^{\frac{1}{2}u}}\right) O_p(n^{(1-\frac{1}{2}u)\alpha_k}) = \mathbf{1}_{\{u>0\}} O_p\left(\frac{n^{(1+\frac{1}{2}u)\alpha_k-\tau u}}{\log(n)^{\frac{1}{2}u}}\right) \quad (285)$$

Using the derivations above is straightforward to see that $III = O_p(II)$.

I therefore conclude that, for $\alpha_k > \max\{\frac{1+\tau}{3}, \frac{\alpha_{max}+4\tau}{6}\}$:

$$T_{zk}^u - \hat{T}_{zk}^u = n^{\frac{1}{2}u} \left[I + II + III \right] \quad (286)$$

$$\begin{aligned} &= n^{\frac{1}{2}u} \left[O_p(n^{(\frac{1}{2}-\frac{1}{2}u)\alpha_k}) + \mathbf{1}_{\{u>0\}} [\bar{O}_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) \right. \\ &\quad \left. + \bar{O}_p(n^{1-2\alpha_k})] O_p(n^{(1-\frac{1}{2}u)\alpha_k}) \right] \end{aligned} \quad (287)$$

$$\begin{aligned} &= n^{\frac{1}{2}u} \left[O_p(n^{(\frac{1}{2}-\frac{1}{2}u)\alpha_k}) + \mathbf{1}_{\{u>0\}} [O_p(\min\{n^{(1-\frac{1}{2}u)\alpha_k}, n^{(1-\frac{1}{2}u)\alpha_k + \frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}\}) \right. \\ &\quad \left. + O_p(\min\{n^{(1-\frac{1}{2}u)\alpha_k}, n^{(1-\frac{1}{2}u)\alpha_k + 1-2\alpha_k}\})] \right] \end{aligned} \quad (288)$$

$$\begin{aligned} &= n^{(1-\frac{1}{2}u)\alpha_k + \frac{1}{2}u} \left[O_p(n^{-\frac{1}{2}\alpha_k}) + \mathbf{1}_{\{u>0\}} [O_p(\min\{1, n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}\}) \right. \\ &\quad \left. + O_p(\min\{1, n^{1-2\alpha_k}\})] \right] \end{aligned} \quad (289)$$

$$= n^{(1-\frac{1}{2}u)\alpha_k + \frac{1}{2}u} \left[O_p(n^{-\frac{1}{2}\alpha_k}) + \mathbf{1}_{\{u>0\}} [\bar{O}_p(n^{\frac{1}{4}\alpha_{max}-\frac{1}{2}\alpha_k}) \bar{O}_p(n^{1-2\alpha_k})] \right] \quad (290)$$

For part (c), with $\alpha_k \leq \max\{\frac{1+\tau}{3}, \frac{\alpha_{max}+4\tau}{6}\}$:

$$T_{zk}^u - \hat{T}_{zk}^u = n^{\frac{1}{2}u} \left[I + II + III \right] \quad (291)$$

$$= n^{\frac{1}{2}u} \left[O_p\left(\frac{n^{(\frac{1}{2}+\frac{1}{2}u)\alpha_k-\tau u}}{\log(n)^{\frac{1}{2}u}}\right) + \mathbf{1}_{\{u>0\}} O_p\left(\frac{n^{(1+\frac{1}{2}u)\alpha_k-\tau u}}{\log(n)^{\frac{1}{2}u}}\right) \right] \quad (292)$$

$$= \frac{n^{(1+\frac{1}{2}u)\alpha_k+(\frac{1}{2}-\tau)u}}{\log(n)^{\frac{1}{2}u}} \left[O_p(n^{-\frac{1}{2}\alpha_k}) + \mathbf{1}_{\{u>0\}} O_p(1) \right] \quad (293)$$

We conclude, combining the above with Lemma 2, that

1. For $\alpha_k > \max\{\frac{1+\tau}{3}, \frac{\alpha_{max}+4\tau}{6}\}$: $\hat{T}_{zk}^u = T_{zk}^u + o_p(T_{zk}^u)$,
2. For $0 < \alpha_k \leq \max\{\frac{1+\tau}{3}, \frac{\alpha_{max}+4\tau}{6}\}$: $\hat{T}_{zk}^u = T_{zk}^u + \mathbf{1}_{\{u>0\}} O_p(T_{zk}^u) + o_p(T_{zk}^u)$,
3. For $\alpha_k = 0$ and $k > r$: $\hat{T}_{zk}^u = O_p(n^{(\frac{1}{2}-\tau)u} \log(n)^{-\frac{1}{2}u})$,

which finishes the proof. \square

D.3 Arbitrage Pricing Theory

Proof of Proposition 1. The proof largely follows the proof of Theorem 3 in Green and Hollifield (1992). Define the set of demeaned portfolios

$$\Xi_n = \{R^p - \mathbb{E}(R^p) : R^p = \sum_{i=1}^n w_i R_i, \sum_{i=1}^n w_i = 1\} \quad (294)$$

and construct the factor-mimicking portfolios by projecting the zero vector and the strongest K factors $k = 1, \dots, K$ onto Ξ_n , such that:

$$F_k = R_{nk}^* - \mathbb{E}(R_{nk}^*) + \xi_{nk}, \quad (295)$$

where $E(\xi_{nk} R_j) = 0$ for $j = 1, \dots, n$ (Noting that R_{n0}^* will be the minimum-variance portfolio with zero betas). For asset j , consider the combination of K factor-mimicking portfolios with the same factor risk:

$$R_{nj}^{K*} = (1 - \sum_k^K \lambda_{jk}) R_{n0}^* + \sum_k^K \lambda_{jk} R_{nk}^*. \quad (296)$$

Let

$$\Pi_{nj}^K = R_j - R_{nj}^{K*} \quad (297)$$

$$= R_j - \mathbb{E}(R_j) + \mathbb{E}(R_j) - (1 - \sum_k^K \lambda_{jk}) R_{n0}^* - \sum_k^K \lambda_{jk} R_{nk}^* \quad (298)$$

$$= c_j^K + \sum_k^K \lambda_{jk} F_k + e_j^K - (1 - \sum_k^K \lambda_{jk}) [R_{n0}^* - \mathbb{E}(R_{n0}^*)] - \sum_k^K \lambda_{jk} [R_{nk}^* - \mathbb{E}(R_{nk}^*)] \quad (299)$$

$$= c_j^K + (1 - \sum_k^K \lambda_{jk}) \xi^{n0} + \sum_k^K \lambda_{jk} \xi_{nk} + e_j^K, \quad (300)$$

with

$$c_j^K = \mathbb{E}(R_j) - \left((1 - \sum_k^K \lambda_{jk}) \mathbb{E}(R_{n0}^*) + \sum_k^K \lambda_{jk} \mathbb{E}(R_{nk}^*) \right). \quad (301)$$

Recalling that W_n denotes the sup-norm on the asset weights w_i , we can invoke the following result by Green and Hollifield (1992).

Theorem (Theorem 1 of Green and Hollifield (1992)). *The efficient portfolio with mean $\mu \neq \nu$ is*

well diversified (i.e. $|w_i| \leq W_n \forall i$) if and only if the return, R^* , on every portfolio with weights that sum to one, satisfies

$$|\mathbb{E}(R^*) - \mathbb{E}(R_z)| \leq \left| \frac{W_n}{\gamma_n} \right| \sum_{i=1}^n |Cov(R^*, R_i)|$$

and the payoff, Π^* , on every hedge position, with weights that sum to zero, satisfies

$$|\mathbb{E}(\Pi^*)| \leq \left| \frac{W_n}{\gamma_n} \right| \sum_{i=1}^n |Cov(\Pi^*, R_i)|,$$

where γ_n is uniformly bounded away from zero by the assumption of no asymptotic arbitrage.

Therefore, if the efficient frontier contains a well-diversified portfolio, this implies that

$$|\mathbb{E}(\Pi_{nj}^K)| \leq \left| \frac{W_n}{\gamma_n} \right| \sum_{i=1}^n |Cov(\Pi_{nj}^K, R_i)|, \quad (302)$$

because Π_j^n is the return on a hedge position with weights summing to zero. By (300), $Cov(\Pi_{nj}^K, R_i) = Cov(e_i^K, e_j^K)$ and thus:

$$|\mathbb{E}(\Pi_{nj}^K)| \leq \left| \frac{W_n}{\gamma_n} \right| \sum_{i=1}^n |Cov(e_i^K, e_j^K)| \quad (303)$$

$$= \left| \frac{W_n}{\gamma_n} \right| \sum_{i=1}^n |Cov(F^w \lambda_i + e_i, F^w \lambda_j + e_j)| \quad (304)$$

$$\leq \left| \frac{W_n}{\gamma_n} \right| \left(\sum_{i=1}^n \sum_{k=K+1}^r \lambda_{ik} \lambda_{jk} + \sum_{i=1}^n |Cov(e_i, e_j)| \right) \quad (305)$$

$$= \left| \frac{W_n}{\gamma_n} \right| \left(\sum_{k=K+1}^r \lambda_{jk} \sum_{i=1}^n \lambda_{ik} + \sum_{i=1}^n |Cov(e_i, e_j)| \right) \quad (306)$$

$$= \left| \frac{W_n}{\gamma_n} \right| \left(\sum_{k=K+1}^r \lambda_{jk} \left[\sum_{i \in \mathcal{A}_k} \lambda_{ik} + \sum_{i \notin \mathcal{A}_k} \lambda_{ik} \right] + \sum_{i=1}^n |Cov(e_i, e_j)| \right) \quad (307)$$

$$= \left| \frac{W_n}{\gamma_n} \right| \left(\sum_{k=K+1}^r \lambda_{jk} [O(n^{\alpha_k}) + O(\sqrt{n})] + \sum_{i=1}^n |Cov(e_i, e_j)| \right) \quad (308)$$

$$\leq \left| \frac{W_n}{\gamma_n} \right| \left(\sum_{k=K+1}^r O(n^{\alpha_k}) + O(\sqrt{n}) + O(\sqrt{n}) \right) \quad (309)$$

$$\leq \left| \frac{W_n}{\gamma_n} \right| \left(O(n^{\alpha_{K+1}}) + O(\sqrt{n}) + O(\sqrt{n}) \right), \quad (310)$$

I therefore conclude that

$$\lim_{n \rightarrow \infty} \left[\mathbb{E}(R_j) - \left((1 - \sum_k^K \lambda_{jk}) \mathbb{E}(R_{n0}^*) + \sum_k^K \lambda_{jk} \mathbb{E}(R_{nk}^*) \right) \right] \quad (311)$$

$$= \lim_{n \rightarrow \infty} |\mathbb{E}(\Pi_{nj}^K)| = \lim_{n \rightarrow \infty} W_n O(\max(n^{\alpha_K+1}, \sqrt{n})) = 0, \quad (312)$$

whenever $W_n = o(\min(n^{-\alpha_K+1}, n^{-\frac{1}{2}}))$. This completes the proof. \square

References

- Ahn, S. C. and A. R. Horenstein (2013). Eigenvalue ratio test for the number of factors. *Econometrica* 81(3), 1203–1227.
- Alessi, L., M. Barigozzi, and M. Capasso (2010). Improved penalization for determining the number of factors in approximate factor models. *Statistics & Probability Letters* 80(23), 1806–1813.
- 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(519), 1182–1198.
- Andrews, D. W. and X. Cheng (2012). Estimation and inference with weak, semi-strong, and strong identification. *Econometrica* 80(5), 2153–2211.
- Antoine, B. and E. Renault (2012). Efficient minimum distance estimation with multiple rates of convergence. *Journal of Econometrics* 170(2), 350–367.
- Bai, J. (2003). Inferential theory for factor models of large dimensions. *Econometrica* 71(1), 135–171.
- Bai, J. and S. Ng (2002). Determining the number of factors in approximate factor models. *Econometrica* 70(1), 191–221.
- Bai, J. and S. Ng (2006a). Confidence intervals for diffusion index forecasts and inference for factor-augmented regressions. *Econometrica* 74(4), 1133–1150.
- Bai, J. and S. Ng (2006b). Determining the number of factors in approximate factor models, errata. Technical report, Columbia University.
- Bernanke, B., J. Boivin, and P. Eliasz (2005). Factor augmented vector autoregressions (fvars) and the analysis of monetary policy. *Quarterly Journal of Economics* 120(1), 387–422.
- Boivin, J. and S. Ng (2006). Are more data always better for factor analysis? *Journal of Econometrics* 132(1), 169–194.
- Cai, T. T., Z. Ma, and Y. Wu (2013). Sparse PCA: Optimal rates and adaptive estimation. *The Annals of Statistics* 41(6), 3074–3110.
- Carvalho, C. M., J. Chang, J. E. Lucas, J. R. Nevins, Q. Wang, and M. West (2008). High-dimensional sparse factor modeling: Applications in gene expression genomics. *Journal of the American Statistical Association* 103(484), 1438–1456.

- Cattell, R. B. (1966). The scree test for the number of factors. *Multivariate behavioral research* 1(2), 245–276.
- Chamberlain, G. (1983). Funds, factors, and diversification in arbitrage pricing models. *Econometrica* 51(5), 1305–1323.
- Chamberlain, G. and M. Rothschild (1983). Arbitrage, factor structure, and mean-variance analysis on large asset markets. *Econometrica* 51(5), 1281–1304.
- Choi, I. (2012). Efficient estimation of factor models. *Econometric Theory* 28(2), 274–308.
- Chudik, A., M. H. Pesaran, and E. Tosetti (2011). Weak and strong cross-section dependence and estimation of large panels. *The Econometrics Journal* 14(1), C45–C90.
- Connor, G. and R. A. Korajczyk (1993). A test for the number of factors in an approximate factor model. *The Journal of Finance* 48(4), 1263–1291.
- Connor, G. and R. A. Korajczyk (1995). The arbitrage pricing theory and multifactor models of asset returns. In V. M. R.A. Jarrow and W. Ziemba (Eds.), *Handbooks in operations research and management science*, Volume 9, Chapter 4, pp. 87–144. Elsevier.
- 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(2), 318–328.
- DeMiguel, V., L. Garlappi, and R. Uppal (2009). Optimal versus naive diversification: How inefficient is the 1/n portfolio strategy? *Review of Financial Studies* 22(5), 1915–1953.
- Dias, F., M. Pinheiro, and A. Rua (2013). Determining the number of global and country-specific factors in the euro area. *Studies in Nonlinear Dynamics and Econometrics* 17(5), 573–617.
- Foerster, A. T., P.-D. G. Sarte, and M. W. Watson (2011). Sectoral versus aggregate shocks: A structural factor analysis of industrial production. *Journal of Political Economy* 119(1), 1–38.
- Forni, M., M. Hallin, M. Lippi, and L. Reichlin (2000). The generalized dynamic-factor model: Identification and estimation. *Review of Economics and statistics* 82(4), 540–554.
- Freyaldenhoven, S. (2017). Sparse factor models. Working paper, Brown University.
- Gabaix, X. (2011). The granular origins of aggregate fluctuations. *Econometrica* 79(3), 733–772.

- Gao, C., C. D. Brown, and B. E. Engelhardt (2013). A latent factor model with a mixture of sparse and dense factors to model gene expression data with confounding effects. *arXiv preprint arXiv:1310.4792*.
- Giannone, D., L. Reichlin, and L. Sala (2006). VARs, common factors and the empirical validation of equilibrium business cycle models. *Journal of Econometrics* 132(1), 257–279.
- Green, R. C. and B. Hollifield (1992). When will mean-variance efficient portfolios be well diversified? *The Journal of Finance* 47(5), 1785–1809.
- Hallin, M. and R. Liska (2007). The generalized dynamic factor model: Determining the number of factors. *Journal of the American Statistical Association* 102(478), 603–617.
- Horn, R. A. and C. R. Johnson (2012). *Matrix analysis*. Cambridge University Press.
- Horvath, M. (1998). Cyclicalities and sectoral linkages: Aggregate fluctuations from independent sectoral shocks. *Review of Economic Dynamics* 1(4), 781–808.
- Huberman, G., S. Kandel, and R. F. Stambaugh (1987). Mimicking portfolios and exact arbitrage pricing. *The Journal of Finance* 42(1), 1–9.
- Kapetanios, G. (2004). A new method for determining the number of factors in factor models with large datasets. Technical report, Department of Economics, Queen Mary, University of London.
- 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(3), 397–409.
- Kleibergen, F. (2009). Tests of risk premia in linear factor models. *Journal of Econometrics* 149(2), 149–173.
- Long, J. B. and C. I. Plosser (1983). Real business cycles. *Journal of political Economy* 91(1), 39–69.
- Moench, E., S. Ng, and S. Potter (2013). Dynamic hierarchical factor models. *Review of Economics and Statistics* 95(5), 1811–1817.
- Moon, H. R. and M. Weidner (2017). Dynamic linear panel regression models with interactive fixed effects. *Econometric Theory* 33, 158–195.
- Onatski, A. (2009). Testing hypotheses about the number of factors in large factor models. *Econometrica* 77(5), 1447–1479.

- Onatski, A. (2010). Determining the number of factors from empirical distribution of eigenvalues. *The Review of Economics and Statistics* 92(4), 1004–1016.
- Onatski, A. (2012). Asymptotics of the principal components estimator of large factor models with weakly influential factors. *Journal of Econometrics* 168(2), 244–258.
- Onatski, A. (2015). Asymptotic analysis of the squared estimation error in misspecified factor models. *Journal of Econometrics* 186(2), 388–406.
- Pati, D., A. Bhattacharya, N. S. Pillai, and D. Dunson (2014). Posterior contraction in sparse bayesian factor models for massive covariance matrices. *The Annals of Statistics* 42(3), 1102–1130.
- Paul, D. and I. M. Johnstone (2012). Augmented sparse principal component analysis for high dimensional data. *arXiv preprint arXiv:1202.1242*.
- Ross, S. A. (1976). The arbitrage theory of capital asset pricing. *Journal of Economic Theory* 13(3), 341–360.
- Shukla, R. and C. Trzcinka (1990). Sequential tests of the arbitrage pricing theory: A comparison of principal components and maximum likelihood factors. *The Journal of Finance* 45(5), 1541–1564.
- Stock, J. H. and M. W. Watson (2002a). Forecasting using principal components from a large number of predictors. *Journal of the American Statistical Association* 97(460), 1167–1179.
- Stock, J. H. and M. W. Watson (2002b). Macroeconomic forecasting using diffusion indexes. *Journal of Business & Economic Statistics* 20(2), 147–162.
- Stock, J. H. and M. W. Watson (2005). Implications of dynamic factor models for var analysis. Technical report, NBER Working Paper 114467.
- Stock, J. H. and M. W. Watson (2006). Forecasting with many predictors. In C. G. G. Elliott and A. Timmermann (Eds.), *Handbook of Economic Forecasting*, Volume 1, Chapter 10, pp. 515–554. Elsevier.
- Stock, J. H. and M. W. Watson (2012). Disentangling the channels of the 2007-09 recession. *Brookings Papers on Economic Activity* 43, 81–156.
- Stock, J. H. and M. W. Watson (2016). Dynamic factor models, factor-augmented vector autoregressions, and structural vector autoregressions in macroeconomics. In J. B. Taylor and H. Uhlig (Eds.), *Handbook of Macroeconomics*, Volume 2, Chapter 8, pp. 415–525. Elsevier.

Trzcinka, C. (1986). On the number of factors in the arbitrage pricing model. *The Journal of Finance* 41(2), 347–368.

Wang, P. (2008). Large dimensional factor models with a multi-level factor structure: Identification, estimation and inference. Working paper, New York University.