



Learning Latent Factors From Diversified Projections and Its Applications to Over-Estimated and Weak Factors

Jianqing Fan & Yuan Liao

To cite this article: Jianqing Fan & Yuan Liao (2022) Learning Latent Factors From Diversified Projections and Its Applications to Over-Estimated and Weak Factors, Journal of the American Statistical Association, 117:538, 909-924, DOI: [10.1080/01621459.2020.1831927](https://doi.org/10.1080/01621459.2020.1831927)

To link to this article: <https://doi.org/10.1080/01621459.2020.1831927>



View supplementary material [↗](#)



Published online: 20 Nov 2020.



Submit your article to this journal [↗](#)



Article views: 1773



View related articles [↗](#)



View Crossmark data [↗](#)



Citing articles: 4 View citing articles [↗](#)



Learning Latent Factors From Diversified Projections and Its Applications to Over-Estimated and Weak Factors

Jianqing Fan^a and Yuan Liao^b

^aDepartment of Operations Research and Financial Engineering, Princeton University, Princeton, NJ; ^bDepartment of Economics, Rutgers University, New Brunswick, NJ

ABSTRACT

Estimations and applications of factor models often rely on the crucial condition that the number of latent factors is consistently estimated, which in turn also requires that factors be relatively strong, data are stationary and weakly serially dependent, and the sample size be fairly large, although in practical applications, one or several of these conditions may fail. In these cases, it is difficult to analyze the eigenvectors of the data matrix. To address this issue, we propose simple estimators of the latent factors using cross-sectional projections of the panel data, by weighted averages with predetermined weights. These weights are chosen to diversify away the idiosyncratic components, resulting in “diversified factors.” Because the projections are conducted cross-sectionally, they are robust to serial conditions, easy to analyze and work even for finite length of time series. We formally prove that this procedure is robust to over-estimating the number of factors, and illustrate it in several applications, including post-selection inference, big data forecasts, large covariance estimation, and factor specification tests. We also recommend several choices for the diversified weights. Supplementary materials for this article are available online.

ARTICLE HISTORY

Received March 2020
Accepted September 2020

KEYWORDS

Factor-augmented regression; Large dimensions; Over-estimating the number of factors; Principal components; Random projections

1. Introduction

Consider the following high-dimensional factor model:

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

where $\mathbf{x}_t = (x_{1t}, \dots, x_{Nt})'$ is an N -dimensional outcome. In addition, the model contains \mathbf{f}_t as r -dimensional latent factors, $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_N)'$ as $N \times r$ matrix of loadings, and $\mathbf{u}_t = (u_{1t}, \dots, u_{Nt})'$ as idiosyncratic terms. Theoretical studies of the model have been crucially depending on the assumption that the number of factors, r , should be consistently estimated. This in turn, requires the factors be relatively strong, data have weak serial dependence, and length of time series T is long. But in practical applications, one or several of these conditions may fail to hold due to weak signal-to-noise ratios and nonstationary or noisy data, making the first r eigenvalues of the sample covariance of $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_T)$ empirically be not so-well separated from the remaining ones.

A promising remedy is to over-estimate the number of factors. But this approach has been quite challenging. Let R be the “working number of factors” that are empirically estimated. When $R > r$, it is often difficult to analyze the behavior of the $(R - r)$ eigenvalues/eigenvectors. As shown in Johnstone and Lu (2009), these eigenvectors can be inconsistent because their eigenvalues are not so “spiked.” This creates challenges to many factor estimators, such as the popular principal components (PC)-estimator (Connor and Korajczyk 1986; Stock and Watson 2002), and therefore brings obstacles to applications when the number of factors is over-estimated. Another difficulty is to handle the serial dependence. As shown by Bai (2003),

the PC-estimator is inconsistent under finite- T in the presence of serial correlations and heteroscedasticity, but many forecast applications using estimated factors favor relatively short time series, due to the concern of nonstationarity.

This article proposes a new method to address issues of over-estimating the number of factors, small T , and strong serial conditions. We propose a simple factor estimator that does not rely on eigenvectors. Let $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_R)$ be a given exogenous (or deterministic) $N \times R$ matrix, where each of its R columns \mathbf{w}_k is an $N \times 1$ vector of “diversified weights,” in the sense that its strength should be approximately equally distributed on most of its components. We propose to estimate \mathbf{f}_t by simply

$$\hat{\mathbf{f}}_t = \frac{1}{N} \mathbf{W}' \mathbf{x}_t,$$

or more precisely, the linear space spanned by $\{\mathbf{f}_t\}_{t=1}^T$ is estimated by that spanned by $\{\hat{\mathbf{f}}_t\}_{t=1}^T$. By substituting (1) into the definition, we have

$$\hat{\mathbf{f}}_t = \underbrace{\left(\frac{1}{N} \mathbf{W}' \mathbf{B} \right)}_{\text{affine transform}} \mathbf{f}_t + \frac{1}{N} \mathbf{W}' \mathbf{u}_t. \quad (2)$$

Thus, $\hat{\mathbf{f}}_t$ (consistently) estimates \mathbf{f}_t up to an $R \times r$ affine transform, with $\mathbf{e}_t := \frac{1}{N} \mathbf{W}' \mathbf{u}_t$ as the estimation error. The assumption that \mathbf{W} should be diversified ensures that as $N \rightarrow \infty$, \mathbf{e}_t is “diversified away” (converging to zero in probability).

We call the new factor estimator as “diversified factors,” which reduces the dimension of \mathbf{x}_t through diversified projections. Because of the clean expansion (2), the mathematics

for theoretical analysis is much simpler than most benchmark estimators. We show that $\hat{\mathbf{f}}_t$ leads to valid inferences in several factor-augmented models so long as $R \geq r$. Therefore, we formally justify that the use of factor models is robust to over-estimating the number of factors. In particular, we admit $r = 0$ but $R \geq 1$ as a special case. That is, the inference is still valid even if there are no common factors present, but we nevertheless take out estimated factors (for insurance). Furthermore, the projection is conducted on cross-sections, so is not sensitive to serial conditions. We study several applications in detail, including the post-selection inference, big data forecasts, high-dimensional covariance estimation, and factor specification tests.

One of the key assumptions imposed is that while \mathbf{W} diversifies away \mathbf{u}_t , we have

$$\text{rank}\left(\frac{1}{N}\mathbf{W}'\mathbf{B}\right) = r,$$

and the r th smallest singular value of $\frac{1}{N}\mathbf{W}'\mathbf{B}$ does not decay too fast. That is, \mathbf{W} should not diversify away the factor components in the time series. This condition *does not* hold if \mathbf{W} has more than $R - r$ columns that are nearly orthogonal to \mathbf{B} . This is another motivation of using over-estimated factors: if random weights are used the probability that more than $R - r$ columns of \mathbf{W} are nearly orthogonal to the space of \mathbf{B} should be very small.

To satisfy the above conditions on the weights, we rely on external information on the factor loadings, and recommend four choices for the weight matrix. The first choice is the individual-specific characteristics. As documented in semiparametric factor models (Park et al. 2009; Connor, Matthias, and Linton 2012; Fan, Liao, and Wang 2016), factor loadings are often driven by observed characteristics. When these variables are available, they can be naturally used as diversified weights. The second choice is based on rolling window estimations. Consider time series forecasts. To pertain the stationarity assumption, we divide the sampling periods into (I) $t = 1, \dots, T_0$ and (II) $t = T_0 + 1, \dots, T_0 + T$, and only use the most recent T observations from period (II) to learn the latent factors for forecasts. Or consider a time series where a structural break occurs at time T_0 , so the most recent period (II) is of major interest. Assume that the loadings are correlated between the two periods, then the PC-estimated loadings from periods (I) would be a good choice of the diversified weights for period (II). For the third recommendation, when the time series is independent of the initial observation, we can use transformations of \mathbf{x}_0 as the weights. The fourth recommended choice is to use columns of the Walsh–Hadamard matrix from the statistical experimental design to form the diversified weights.

The idea of approximating factors by weighted averages of observations has been applied previously in the literature. In the asset pricing literature, factors are created by weighted averages of a large number of asset returns. There, the weights are also predetermined, adapted to the filtration up to the last observation time. In the common correlated effects (CCE) literature (Pesaran 2006; Chudik, Pesaran, and Tosetti 2011), factors are created using a set of random weights to estimate the effect of observables. There, R equals the dimensions of additionally observed regressors and the outcome variable, and certain rank

conditions about the regressors are required. In the same setting, Westerlund and Urbain (2015) and Karabiyik, Urbain, and Westerlund (2019) compared the cross-sectional average and the PC estimators, and also showed the validity of using $R > r$ number of cross-sectional averages. Moreover, Barigozzi and Cho (2018) proposed a different method to address the issue of over-estimating factors. One of our recommended weights is inspired by their approach. Moon and Weidner (2015) studied the problem in a panel data framework and showed that the inference about the parameter of interest is robust to over-estimating the number of factors. Finally, there is a large literature on estimating the number of factors. See Bai and Ng (2002), Hallin and Liška (2007), Ahn and Horenstein (2013), and Li, Li, and Shi (2017).

The rest of the article is organized as follows. Section 2 explains the key ideas and intuitions in details. Section 3 presents several applications of the diversified factors. Section 4 recommends several choices of the weight matrix. Section 5 conducts extensive simulation studies using various models. All technical proofs are presented in the supplementary appendix.

We use the following notation. For a matrix \mathbf{A} , we use $\lambda_{\min}(\mathbf{A})$ and $\lambda_{\max}(\mathbf{A})$ to denote its smallest and largest eigenvalues. We define the Frobenius norm $\|\mathbf{A}\|_F = \sqrt{\text{tr}(\mathbf{A}'\mathbf{A})}$ and the operator norm $\|\mathbf{A}\| = \sqrt{\lambda_{\max}(\mathbf{A}'\mathbf{A})}$. In addition, define projection matrices $\mathbf{M}_\mathbf{A} = \mathbf{I} - \mathbf{P}_\mathbf{A}$ and $\mathbf{P}_\mathbf{A} = \mathbf{A}(\mathbf{A}'\mathbf{A})^{-1}\mathbf{A}$ when $\mathbf{A}'\mathbf{A}$ is invertible. Finally, for two (random) sequences a_T and b_T , we write $a_T \ll b_T$ (or $b_T \gg a_T$) if $a_T = o_P(b_T)$.

2. Factor Estimation Using Diversified Projections

2.1. The Estimator

Let $R \geq r$ be a predetermined bounded integer that does not grow with N , which we call “the working number of factors.” As in practice we do not know the true number of factors r , we often take a slightly large R so that $R \geq r$ is likely to hold. Let $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_R)$ be a user-specified $N \times R$ matrix, either deterministic or random but independent of the σ -algebra generated by $\{\mathbf{u}_t : t = 1, 2, \dots\}$. Each of its R columns $\mathbf{w}_k = (w_{k,1}, \dots, w_{k,N})'$ ($k \leq R$) is an $N \times 1$ vector satisfying the following:

Assumption 2.1 (Diversified weights). There are constants $0 < c < C$, so that (almost surely if \mathbf{W} is random) as $N \rightarrow \infty$,

- (i) $\max_{i \leq N} |w_{k,i}| < C$.
- (ii) The $R \times R$ matrix $\frac{1}{N}\mathbf{W}'\mathbf{W}$ satisfies $\lambda_{\min}(\frac{1}{N}\mathbf{W}'\mathbf{W}) > c$.
- (iii) \mathbf{W} is independent of $\{\mathbf{u}_t : t \leq T\}$.

Construct a factor estimator as an $R \times 1$ vector at each time t :

$$\hat{\mathbf{f}}_t := \frac{1}{N}\mathbf{W}'\mathbf{x}_t.$$

In financial economics applications where \mathbf{x}_t is a vector of asset returns, then each component of $\hat{\mathbf{f}}_t$ is essentially a diversified portfolio return at time t due to its linear form. The behavior of $\hat{\mathbf{f}}_t$ is strikingly simple and clean. Define an $R \times r$ matrix

$$\mathbf{H} := \frac{1}{N}\mathbf{W}'\mathbf{B}.$$

Then, it follows from the definition and (1) that

$$\widehat{\mathbf{f}}_t = \mathbf{H}\mathbf{f}_t + \frac{1}{N}\mathbf{W}'\mathbf{u}_t. \quad (3)$$

Therefore, $\widehat{\mathbf{f}}_t$ estimates an affine transformation of \mathbf{f}_t , where \mathbf{H} is the $R \times r$ transformation matrix. The estimation error equals the diversified idiosyncratic noise $\frac{1}{N}\mathbf{w}'_k\mathbf{u}_t = \frac{1}{N}\sum_{i=1}^N w_{k,i}u_{it}$ for each $k \leq R$. When (u_{1t}, \dots, u_{Nt}) are cross-sectionally weakly dependent, [Assumption 2.1](#) ensures that $\frac{1}{N}\mathbf{w}'_k\mathbf{u}_t$ admits a cross-sectional central limit theorem. For instance, in the special case of cross-sectional independence, it is straightforward to verify the Lindeberg's condition under [Assumption 2.1](#), and therefore as $N \rightarrow \infty$,

$$\frac{1}{\sqrt{N}}\mathbf{W}'\mathbf{u}_t \xrightarrow{d} \mathcal{N}(0, \mathbf{V}), \quad (4)$$

where $\mathbf{V} = \lim_{N \rightarrow \infty} \frac{1}{N}\mathbf{W}'\text{var}(\mathbf{u}_t)\mathbf{W}$ which is assumed to exist.

The convergence (4) shows that $\sqrt{N}(\widehat{\mathbf{f}}_t - \mathbf{H}\mathbf{f}_t)$ is asymptotically normal for each $t \leq T$. Importantly, it holds regardless of whether $T \rightarrow \infty$, $R = r$, or not. It requires only that $N \rightarrow \infty$ and that the weights should be chosen to satisfy [Assumption 2.1](#). This fact is particularly useful for analyzing short time series.

In addition, the factor components should not be diversified away. This gives rise to the following condition on the transformation matrix. Let $\nu_{\min}(\mathbf{H})$ and $\nu_{\max}(\mathbf{H})$, respectively, denote the minimum and maximum *nonzero* singular values of \mathbf{H} .

Assumption 2.2. Suppose $R \geq r$. Almost surely

- (i) $\text{rank}(\mathbf{H}) = r$.
- (ii) There is $C > 0$,

$$\nu_{\min}^2(\mathbf{H}) \gg \frac{1}{N}, \quad \nu_{\max}(\mathbf{H}) \leq C\nu_{\min}(\mathbf{H}).$$

[Assumption 2.2](#) requires that \mathbf{W} have at least r columns that are not orthogonal to \mathbf{B} so that \mathbf{B} is not diversified away. This is the key assumption, but is not stringent in the context of over-estimating factors. In the current setting, the factor strength is measured by $\nu_{\min}(\mathbf{H})$, which is required not to decay very fast by condition (ii). This quantity determines the rate of convergence in recovering the space spanned by the factors.

Given $\widehat{\mathbf{f}}_t$, it is straightforward to estimate the loading matrix by using the least squares:

$$\widehat{\mathbf{B}} = (\widehat{\mathbf{b}}_1, \dots, \widehat{\mathbf{b}}_N)' = \sum_{t=1}^T \mathbf{x}_t \widehat{\mathbf{f}}_t' \left(\sum_{t=1}^T \widehat{\mathbf{f}}_t \widehat{\mathbf{f}}_t' \right)^{-1}.$$

We show that the $R \times R$ matrix $\frac{1}{T}\sum_{t=1}^T \widehat{\mathbf{f}}_t \widehat{\mathbf{f}}_t'$ is nonsingular with probability approaching one even when $R > r$. So $\widehat{\mathbf{B}}$ is well defined. Finally, \mathbf{u}_t can be estimated by

$$\widehat{\mathbf{u}}_t = (\widehat{u}_{1t}, \dots, \widehat{u}_{Nt}) = \mathbf{x}_t - \widehat{\mathbf{B}}\widehat{\mathbf{f}}_t. \quad (5)$$

Just like the PC-estimator, the diversified projection can estimate dynamic factor models by treating dynamic factors as static factors. In addition, it is straightforward to extend the model to allowing time-varying factor loadings, by time-domain local smoothing before applying the diversified projection. While these extensions are straightforward, here we focus on static and time invariant models.

2.2. Over-Estimating the Number of Factors

The consistent estimation for the number of factors r often requires strong conditions that may be violated in finite sample. An advantage of the diversified factors is being robust to over-estimating the number of factors in many inference problems.

We start with a heuristic discussion of the main issue in this subsection. Recall that $\mathbf{H} = \frac{1}{N}\mathbf{W}'\mathbf{B}$ is the $R \times r$ matrix, which is no longer a square matrix when $R > r$. In this case $\widehat{\mathbf{B}}$ is essentially estimating $\mathbf{B}\mathbf{H}^+$, with the $r \times R$ transformation matrix \mathbf{H}^+ being the Moore–Penrose generalized inverse of \mathbf{H} , defined as follows. Suppose \mathbf{H}' has the following singular value decomposition:

$$\mathbf{H}' = \mathbf{U}_H(\mathbf{D}_H, 0)\mathbf{E}_H', \quad r \times R,$$

where 0 in the above singular value matrix is present whenever $R > r$, and \mathbf{D}_H is an $r \times r$ diagonal matrix of the nonzero singular values. Then \mathbf{H}^+ is an $r \times R$ matrix:

$$\mathbf{H}^+ = \mathbf{U}_H(\mathbf{D}_H^{-1}, 0)\mathbf{E}_H'.$$

It is straightforward to verify that $\mathbf{H}^+\mathbf{H} = \mathbf{I}_r$ holds and that for estimating the common component $\mathbf{B}\mathbf{f}_t$ using over-estimated number of factors, we have

$$\widehat{\mathbf{B}}\widehat{\mathbf{f}}_t = \mathbf{B}\mathbf{H}^+\mathbf{H}\mathbf{f}_t + o_P(1) = \mathbf{B}\mathbf{f}_t + o_P(1), \quad (6)$$

where $o_P(1)$ in the above approximation can be made uniformly across elements.

However, a key challenge of formalizing the intuition behind (6) is to analyze the invertibility of the gram matrix $\frac{1}{T}\sum_{t=1}^T \widehat{\mathbf{f}}_t \widehat{\mathbf{f}}_t'$, which appears in the definition of $\widehat{\mathbf{B}}$. It is also a key ingredient in most applications of factor-augmented models wherever the estimated factors are used as regressors. Define

$$\widehat{\mathbf{S}}_f = \frac{1}{T}\sum_{t=1}^T \widehat{\mathbf{f}}_t \widehat{\mathbf{f}}_t', \quad \mathbf{S}_f = \mathbf{H}\frac{1}{T}\sum_{t=1}^T \mathbf{f}_t \mathbf{f}_t' \mathbf{H}',$$

where \mathbf{S}_f is the population analogue of $\widehat{\mathbf{S}}_f$. The following three bounds when $R > r$, proved in Proposition A.1, play a fundamental role in the asymptotic analysis throughout the article:

- (i) With probability approaching one, $\widehat{\mathbf{S}}_f$ is invertible, but its eigenvalues may decay quickly so that

$$\|\widehat{\mathbf{S}}_f^{-1}\| = O_P(N). \quad (7)$$

On the other hand, \mathbf{S}_f is degenerate when $R > r$, whose rank equals r . Also note that we still have $\|\widehat{\mathbf{S}}_f^{-1}\| = O_P(1)$ when $R = r$ holds.

- (ii) Even if $R > r$, $\|\mathbf{H}'\widehat{\mathbf{S}}_f^{-1}\|$ is much smaller:

$$\|\mathbf{H}'\widehat{\mathbf{S}}_f^{-1}\| = O_P\left(\sqrt{\frac{\max\{N, T\}}{T}}\right).$$

- (iii) When $R > r$, $\|\widehat{\mathbf{S}}_f^{-1} - \mathbf{S}_f^+\| \neq o_P(1)$ but we have

$$\|\mathbf{H}'(\widehat{\mathbf{S}}_f^{-1} - \mathbf{S}_f^+)\mathbf{H}\| = O_P\left(\frac{1}{T} + \frac{1}{N}\right).$$

Therefore, $\widehat{\mathbf{S}}_f$ is invertible, and when weighted by the transformation matrix \mathbf{H}' , its inverse is well behaved and fast converges to the generalized inverse of \mathbf{S}_f , even though \mathbf{S}_f is singular when $R > r$. It is sufficient to consider $\mathbf{H}\widehat{\mathbf{S}}_f^{-1}$ in most factor-augmented inference problems, because in regression models $\widehat{\mathbf{S}}_f^{-1}$ often appears in the projection matrix $\mathbf{P}_{\widehat{\mathbf{F}}} = \widehat{\mathbf{F}}(\widehat{\mathbf{F}}'\widehat{\mathbf{F}})^{-1}\widehat{\mathbf{F}}'$ through $\mathbf{H}\widehat{\mathbf{S}}_f^{-1}$ asymptotically, where $\widehat{\mathbf{F}} := (\widehat{\mathbf{f}}_1, \dots, \widehat{\mathbf{f}}_T)'$ and $\mathbf{F} := (\mathbf{f}_1, \dots, \mathbf{f}_T)'$ denote the estimated and true factor matrices.

Remark 2.1. In the CCE literature (e.g., Pesaran 2006; Chudik, Pesaran, and Tosetti 2011), it has also been claimed that estimating the factors using cross-sectional averages does not require consistently estimating the number of factors. While the claim is true, its proof is not straightforward as $\|\widehat{\mathbf{S}}_f^{-1} - \mathbf{S}_f^{-1}\| \neq O_p(1)\|\widehat{\mathbf{S}}_f - \mathbf{S}_f\|$ when $R > r$. Also see Karabiyik, Reese, and Westerlund (2017) and Karabiyik, Urbain, and Westerlund (2019) for more discussions on the related issue. Our method therefore also potentially contributes to this literature as an alternative rigorous approach.

2.3. Estimating the Factor Space

Throughout the article, the loading matrix \mathbf{B} can be either deterministic or random. When they are random, it is assumed that it is independent of \mathbf{u}_t , and all the expectations throughout the article is taken conditionally on \mathbf{B} .

We make the following conditions.

Assumption 2.3.

- (i) $\{(\mathbf{f}_t, \mathbf{u}_t) : t \leq T\}$ is a stationary process, satisfying $\mathbb{E}(\mathbf{u}_t | \mathbf{f}_t) = 0$.
- (ii) There are constants $c, C > 0$, so that $\max_{i \leq N} \|\mathbf{b}_i\| < C$, and almost surely

$$c < \lambda_{\min} \left(\frac{1}{T} \sum_{t=1}^T \mathbf{f}_t \mathbf{f}_t' \right) \leq \lambda_{\max} \left(\frac{1}{T} \sum_{t=1}^T \mathbf{f}_t \mathbf{f}_t' \right) < C.$$

Assumption 2.4 (Weak dependence). There is a constant $C > 0$,

- (i) $\max_{j,i \leq N} \frac{1}{NT} \sum_{q,v \leq N} \sum_{t,s \leq T} |\text{cov}(u_{it}u_{qt}, u_{js}u_{vs} | \mathbf{F})| < C$ almost surely in \mathbf{F} ,
- (ii) $\frac{1}{T} \sum_{s=1}^T \sum_{t=1}^T \mathbb{E} \|\mathbf{f}_t\| \|\mathbf{f}_s\| \|\mathbb{E}(\mathbf{u}_t \mathbf{u}_s' | \mathbf{F})\| < C$ and $\mathbb{E} \|\mathbb{E}(\mathbf{u}_t \mathbf{u}_t' | \mathbf{F})\| < C$.

Theorem 2.1. Suppose Assumptions 2.1–2.4 hold. Also $N \rightarrow \infty$ and T is either finite or grows. Then for all bounded $R \geq r$,

$$\|\mathbf{P}_{\widehat{\mathbf{F}}} \mathbf{P}_{\mathbf{F}} - \mathbf{P}_{\mathbf{F}}\| = O_p \left(\frac{1}{\sqrt{N}} v_{\min}^{-1}(\mathbf{H}) \right), \quad (8)$$

$$\|\mathbf{P}_{\widehat{\mathbf{F}}} - \mathbf{P}_{\mathbf{F}}\| = O_p \left(\frac{1}{\sqrt{N}} v_{\min}^{-1}(\mathbf{H}) \right), \quad (9)$$

where $\mathbf{M} = (\mathbf{H}\mathbf{H}')^+ \mathbf{H}$ is an $R \times r$ matrix¹.

Equation (8) shows that when $R \geq r$, the linear space spanned by $\widehat{\mathbf{F}}$ asymptotically covers the linear space spanned by

¹We show in the proof that $(\mathbf{M}'\widehat{\mathbf{F}}\mathbf{M})$ and $\widehat{\mathbf{F}}'\widehat{\mathbf{F}}$ are both invertible with probability approaching one. So $\mathbf{P}_{\widehat{\mathbf{F}}}$ and $\mathbf{P}_{\mathbf{F}}$ are well defined asymptotically.

\mathbf{F} . To understand the intuition, note that (8) implies $\mathbf{P}_{\widehat{\mathbf{F}}} \mathbf{P}_{\mathbf{F}} \mathbf{Y} \approx \mathbf{P}_{\mathbf{F}} \mathbf{Y}$ for an arbitrary random matrix \mathbf{Y} . Meanwhile, if we heuristically regard $\mathbf{P}_{\mathbf{F}}$ and $\mathbf{P}_{\widehat{\mathbf{F}}}$ as conditional expectations given \mathbf{F} and $\widehat{\mathbf{F}}$, then approximately,

$$\mathbb{E} \left(\mathbb{E}(\mathbf{Y} | \mathbf{F}) \middle| \widehat{\mathbf{F}} \right) \approx \mathbb{E}(\mathbf{Y} | \mathbf{F}). \quad (10)$$

Let $\text{span}(\mathbf{A})$ denote the linear space spanned by the columns of \mathbf{A} . The approximation (10) is well known to be the “tower property,” which heuristically means $\text{span}(\mathbf{F}) \subseteq \text{span}(\widehat{\mathbf{F}})$.

Equation (9) shows that a particular subspace of $\text{span}(\widehat{\mathbf{F}})$ is consistent for $\text{span}(\mathbf{F})$. In the special case $R = r$, we have $\mathbf{P}_{\widehat{\mathbf{F}}} = \mathbf{P}_{\mathbf{F}}$ since \mathbf{M} in (9) is invertible. It then reduces to the usual space consistency. Importantly, we allow T to be finite.

To gain more insights of these results, let us compare with the usual methods based on estimating the number of factors, for example, the eigenvalue-ratio method of Ahn and Horenstein (2013). There are two key quantities in this comparison: the strength of the spiked eigenvalues of $\mathbf{S}_x := \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t'$, and the largest eigenvalue of $\mathbf{S}_u := \frac{1}{T} \sum_{t=1}^T \mathbf{u}_t \mathbf{u}_t'$.

We consider a setting where we can easily quantify the signal-to-noise ratio, as given in the following example.

Example 2.1. This example presents a *pervasive factor model* that satisfies Assumption 2.2. Suppose each individual loading satisfies $\mathbf{b}_i = v_N \boldsymbol{\lambda}_i$ for some sequence $v_N \asymp N^{-(1-\alpha)/2}$ and $\alpha \in (0, 1]$, where $\{\boldsymbol{\lambda}_i : i \leq N\}$ is a sequence of $r \times 1$ vectors such that:

- (i) $\frac{1}{N} \sum_{i=1}^N \boldsymbol{\lambda}_i \boldsymbol{\lambda}_i' \rightarrow \mathbf{C}$ (or converges in probability if $\boldsymbol{\lambda}_i$ is random) for some positive definite matrix \mathbf{C} ;
- (ii) $v_{\min}(\frac{1}{N} \mathbf{W}' \boldsymbol{\Lambda})$ is bounded away from zero, where $\boldsymbol{\Lambda} = (\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_N)'$.

Then Assumption 2.2 holds for $v_{\min}(\mathbf{H}) \asymp v_N$ and any $\alpha \in (0, 1]$. It is straightforward to verify that the r th spiked eigenvalue satisfies

$$\lambda_r(\mathbf{S}_x) \asymp N^\alpha, \quad \alpha \in (0, 1].$$

Theorem 2.1 then shows that $\|\mathbf{P}_{\widehat{\mathbf{F}}} - \mathbf{P}_{\mathbf{F}}\| = o_p(1)$ for any $\alpha > 0$. To verify condition (ii) in this example, consider a “characteristic based” model described in Section 4, where the baseline loading can be decomposed as $\boldsymbol{\lambda}_i = \mathbf{g}(\mathbf{z}_i) + \boldsymbol{\gamma}_i$, with $\mathbb{E}(\boldsymbol{\gamma}_i | \mathbf{z}_i) = 0$; $\mathbf{g}(\mathbf{z}_i)$ is a nonparametric function of some observable characteristic \mathbf{z}_i , and $\boldsymbol{\gamma}_i$ is the loading components that is orthogonal to the characteristic effects. (See more detailed motivations of this model in Section 4.1) Now take $\mathbf{w}_i = \phi(\mathbf{z}_i)$ as an R -dimensional transformation using R predetermined basis functions ϕ . Then a sufficient condition for (ii) is that $v_{\min}(\mathbf{A})$ is bounded away from zero, where $\mathbf{A} := \frac{1}{N} \sum_{i=1}^N \phi(\mathbf{z}_i) \mathbf{g}(\mathbf{z}_i)'$. In addition, Assumption 2.2(ii) holds as long as $v_{\max}(\mathbf{A}) < C$.²

²Suppose $\{\boldsymbol{\gamma}_i : i \leq N\}$ are cross-sectionally conditionally weakly dependent given $\mathbf{Z} = (\mathbf{z}_i : i \leq N)$. Then $\|\frac{1}{N} \sum_i \phi(\mathbf{z}_i) \boldsymbol{\gamma}_i\|^2 = O_p(X) \frac{1}{N} \sum_i \|\phi(\mathbf{z}_i)\|^2 = o_p(1)$ given the assumption that $X := \max_i \frac{1}{N} \sum_j \|\mathbb{E}[\boldsymbol{\gamma}_i \boldsymbol{\gamma}_j' | \mathbf{Z}]\| = o_p(1)$, which holds if $\boldsymbol{\gamma}_i$ are conditionally weakly correlated. Then $v_{\min}(\frac{1}{N} \mathbf{W}' \boldsymbol{\Lambda}) \geq v_{\min}(\mathbf{A}) - \|\frac{1}{N} \sum_i \phi(\mathbf{z}_i) \boldsymbol{\gamma}_i\| \geq c - o_p(1)$. In addition, $v_{\max}(\mathbf{H}) = v_{\max}(\frac{1}{N} \mathbf{W}' \boldsymbol{\Lambda}) v_N \leq [v_{\max}(\mathbf{A}) + o_p(1)] v_N \leq C v_N \leq C c^{-1} v_{\min}(\mathbf{A}) v_N \leq C c^{-1} [v_{\min}(\frac{1}{N} \mathbf{W}' \boldsymbol{\Lambda}) + o_p(1)] v_N \leq 2 C c^{-1} v_{\min}(\mathbf{H})$.

The key implication of [Example 2.1](#) is that the strength of the spiked eigenvalues can grow at an arbitrarily slow polynomial rate in N , and T is allowed to be finite. In applications where $T \rightarrow \infty$ is required, the growth requirement of T can be very mild. For instance, as we shall show in the high-dimensional factor-augmented regression ([Section 3.2](#)), it is only required that $\log^2 N = o(T)$ if the number of “important” control variables (corresponding to nonzero coefficients) is finite. The relative flexibility on the growth of T is achieved thanks to the fact that the diversified projection does not demand strong eigenvalues of the population covariance matrix.

Now let us revisit the conditions required by the eigenvalue-ratio method by Ahn and Horenstein ([2013](#)). If \mathbf{u}_t is sub-Gaussian, under weak dependence conditions,

$$\lambda_{\max}(\mathbf{S}_u) = O_P\left(\frac{\max\{T, N\}}{T}\right).$$

The selection consistency requires $\lambda_r(\mathbf{S}_x) \gg \lambda_{\max}(\mathbf{S}_u)$, which in this context, becomes $T \gg N^{1-\alpha}$. In the case that the spiked eigenvalues are not so strong ($\alpha < 0.5$), it requires a considerably longer time series to override the effect of the idiosyncratic noise.

2.4. Summary of Advantages

Below we summarize key advantages of the use of diversified projection.

1. It is computationally and mathematically simple.
2. When the true number of factors is over estimated ($R \geq r$), inferences about transformation invariant parameters are still asymptotically valid. This leads to important implications on factor-augmented inferences and out-of sample forecasts.
3. It admits an interesting special case, where $r = 0$ and $R \geq 1$. That is, \mathbf{x}_t is in fact weakly dependent, but we nevertheless estimate “factors.” The resulting inference is still asymptotically valid in this case. We shall formally prove this in the high-dimensional factor-augmented inference in the next section. This shows that extracting estimated factors is a robust inference procedure.
4. As the diversified projections are applied cross-sectionally, some conditions that are needed for the PC-estimator can be weakened. For instance, the space spanned by the latent factors can be consistently estimated even if T is finite. It is also a good choice under weak signal-to-noise ratios where the consistent selection of the number of factors is hard to achieve.
5. After applying the diversified projection to \mathbf{x}_t to reduce to a lower dimensional space, one can continue to employ the PCA on $\hat{\mathbf{f}}_t$ to estimate the factor space and the number of factors. This becomes a low-dimensional PCA problem, and potentially much easier than benchmark methods dealing with large dimensional datasets.

3. Applications

We present several applications of the new diversified factors. Besides those imposed in [Section 2](#), additional assumptions are required in each of these examples. These assumptions are

application-specific and are required even if the oracle number of factors were available.

3.1. Forecasts Using Augmented Factor Regression

Consider forecasting time series using a large panel of augmented factor regression

$$y_{t+h} = \alpha' \mathbf{f}_t + \beta' \mathbf{g}_t + \varepsilon_{t+h}, \quad t = 1, \dots, T$$

$$\mathbf{x}_t = \mathbf{B} \mathbf{f}_t + \mathbf{u}_t$$

with observed data $\{(y_t, \mathbf{x}_t) : t \leq T\}$. Here, $h \geq 0$ is the lead time and \mathbf{g}_t is a vector of observed predictors including lagged outcome variables. The goal is the mean forecast:

$$y_{T+h|T} := \alpha' \mathbf{f}_T + \beta' \mathbf{g}_T := \delta' \mathbf{z}_T,$$

where $\mathbf{z}_t = (\mathbf{f}_t' \mathbf{H}', \mathbf{g}_t')'$ and $\delta' = (\alpha' \mathbf{H}', \beta')$. The prediction also depends on unobservable factors \mathbf{f}_t whose information is contained in a high-dimensional panel of data. This model has been studied extensively in the literature (see, e.g., [Stock and Watson 2002](#); [Bai and Ng 2006](#); [Ludvigson and Ng 2007](#)), where \mathbf{f}_T is replaced by a consistent estimator. Once estimated factors $\hat{\mathbf{f}}_t$ is obtained, the forecast of $y_{T+h|T}$ is straightforward:

$$\hat{y}_{T+h|T} = \hat{\delta}' \hat{\mathbf{z}}_T, \quad \hat{\delta} = \left(\sum_{t=1}^{T-h} \hat{\mathbf{z}}_t \hat{\mathbf{z}}_t' \right)^{-1} \sum_{t=1}^{T-h} \hat{\mathbf{z}}_t y_{t+h},$$

where $\hat{\mathbf{z}}_t = (\hat{\mathbf{f}}_t' \mathbf{H}', \mathbf{g}_t')'$ denotes the estimated regressors. Note that $(\sum_{t=1}^{T-h} \hat{\mathbf{z}}_t \hat{\mathbf{z}}_t')^{-1}$ is well defined even if $R > r$ with high probability. This follows from the invertibility of $\hat{\mathbf{F}}' \mathbf{M}_G \hat{\mathbf{F}}$, a claim to be proved (the definition of \mathbf{G} is clear below, and the notation \mathbf{M}_G is defined in [Section 1](#)).

Our study is motivated by two important yet unsolved issues. First, the study of prediction rates has been crucially relying on the assumption that the number of latent factors is correctly estimated. Second, the time series that are being studied are often relatively short, to preserve the stationarity. As we explained in [Section 2](#), this leads to strong conditions on the strength of factors of using the PC estimator.

We show below that by allowing $R > r$, the diversified projection does not require a consistent estimator of the number of factors. In addition to the assumptions in [Section 2](#), we impose the following conditions on the forecast equation for y_{t+h} . Let \mathbf{G} be the matrix of $\{\mathbf{g}_t : t \leq T - h\}$.

Assumption 3.1.

- (i) $\{\varepsilon_t, \mathbf{f}_t, \mathbf{g}_t, \mathbf{u}_t : t = 1, \dots, T + h\}$ is stationary with $\mathbb{E}(\mathbf{u}_t | \mathbf{f}_t, \mathbf{g}_t) = 0$ and $\mathbb{E}(\varepsilon_t | \mathbf{f}_t, \mathbf{g}_t, \mathbf{u}_t, \mathbf{W}) = 0$.
- (ii) Weak dependence: there is $C > 0$, $\max_{s \leq T} \sum_{t \leq T} |\mathbb{E}(\varepsilon_t \varepsilon_s | \mathbf{F}, \mathbf{G}, \mathbf{W})| < C$ almost surely.
- (iii) Moment bounds: there are $c, C > 0$, $\lambda_{\min}(\frac{1}{T} \mathbf{F}' \mathbf{M}_G \mathbf{F}) > c$, $\lambda_{\min}(\frac{1}{T} \mathbf{G}' \mathbf{M}_{\mathbf{F}\mathbf{H}} \mathbf{G}) > c$, and $c < \lambda_{\min}(\frac{1}{T} \mathbf{G}' \mathbf{G}) \leq \lambda_{\max}(\frac{1}{T} \mathbf{G}' \mathbf{G}) < C$.

Our theory *does not* follow from the standard theory of linear models of [Bai and Ng \(2006\)](#). A new technical phenomenon arises when $R > r$ due to the degeneracy of the gram matrices.

Define $\widehat{\mathbf{Z}} = (\widehat{\mathbf{z}}_1, \dots, \widehat{\mathbf{z}}_{T-h})'$, $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_{T-h})'$ and consider two gram matrices

$$\widehat{\mathbf{Z}}'\widehat{\mathbf{Z}} = \begin{pmatrix} \widehat{\mathbf{F}}'\widehat{\mathbf{F}} & \widehat{\mathbf{F}}'\mathbf{G} \\ \mathbf{G}'\widehat{\mathbf{F}} & \mathbf{G}'\mathbf{G} \end{pmatrix}, \quad \mathbf{Z}'\mathbf{Z} = \begin{pmatrix} \mathbf{H}\mathbf{F}'\mathbf{F}\mathbf{H}' & \mathbf{H}\mathbf{F}'\mathbf{G} \\ \mathbf{G}'\mathbf{F}\mathbf{H}' & \mathbf{G}'\mathbf{G} \end{pmatrix}.$$

The linear regression theory crucially depends on the inverse of $\widehat{\mathbf{Z}}'\widehat{\mathbf{Z}}$, whose population version $\mathbf{Z}'\mathbf{Z}$, in this context, becomes degenerate when $R > r$. The full rank matrix $\frac{1}{T}\widehat{\mathbf{F}}'\mathbf{M}_\mathbf{G}\widehat{\mathbf{F}}$ converges to a degenerate matrix $\mathbf{H}\frac{1}{T}\mathbf{F}'\mathbf{M}_\mathbf{G}\mathbf{F}\mathbf{H}'$, and therefore in general

$$\left\| \left(\frac{1}{T}\widehat{\mathbf{Z}}'\widehat{\mathbf{Z}} \right)^{-1} - \left(\frac{1}{T}\mathbf{Z}'\mathbf{Z} \right)^+ \right\| \neq o_p(1).$$

We develop a new theory that takes advantage of \mathbf{H} , which allows to establish the three claims in Section 2.2. They imply that the convergence holds when weighted by $\widetilde{\mathbf{H}}$:

$$\left\| \widetilde{\mathbf{H}}' \left(\left(\frac{1}{T}\widehat{\mathbf{Z}}'\widehat{\mathbf{Z}} \right)^{-1} - \left(\frac{1}{T}\mathbf{Z}'\mathbf{Z} \right)^+ \right) \widetilde{\mathbf{H}} \right\| = O_p \left(\frac{1}{T} + \frac{1}{N} \right),$$

where $\widetilde{\mathbf{H}} = \begin{pmatrix} \mathbf{H} \\ \mathbf{I} \end{pmatrix}$.

The weighted convergence is sufficient to derive the prediction rate of $\widehat{\mathbf{y}}_{T+h|T}$.

Theorem 3.1. Suppose Assumptions 2.1–2.4 and 3.1 hold. As $T, N \rightarrow \infty$, h is bounded, and for all bounded $R \geq r$,

$$\widehat{\mathbf{y}}_{T+h|T} - \mathbf{y}_{T+h|T} = O_p \left(\frac{1}{\sqrt{T}} + \frac{1}{v_{\min}\sqrt{N}} \right).$$

3.2. High-Dimensional Inference in Factor Augmented Models

3.2.1. Factor-Augmented Post-Selection Inference

Consider a high-dimensional regression model

$$y_t = \beta \mathbf{g}_t + \mathbf{v}'\mathbf{x}_t + \eta_t, \quad (11)$$

$$\mathbf{g}_t = \theta'\mathbf{x}_t + \boldsymbol{\varepsilon}_{g,t}, \quad (12)$$

where \mathbf{g}_t is a treatment variable whose effect β is the main interest. The model contains high-dimensional control variables $\mathbf{x}_t = (x_{1t}, \dots, x_{Nt})'$ that determine both the outcome and treatment variables. Having many control variables creates challenges for statistical inferences, as such, we assume that (\mathbf{v}, θ) are sparse vectors. Belloni, Chernozhukov, and Hansen (2014) proposed to make inference using Robinson's (1988) residual-regression, by first selecting among the high-dimensional controls in both the y_t and \mathbf{g}_t equations.

Often, the control variables are strongly correlated due to the presence of confounding factors

$$\mathbf{x}_t = \mathbf{B}\mathbf{f}_t + \mathbf{u}_t. \quad (13)$$

This invalidates the conditions of using penalized regressions to directly select among \mathbf{x}_t . Instead, if we substitute (13) to (11), we reach factor-augmented regression model:

$$y_t = \boldsymbol{\alpha}'_y \mathbf{f}_t + \boldsymbol{\gamma}'\mathbf{u}_t + \boldsymbol{\varepsilon}_{y,t}, \quad (14)$$

$$\mathbf{g}_t = \boldsymbol{\alpha}'_g \mathbf{f}_t + \boldsymbol{\theta}'\mathbf{u}_t + \boldsymbol{\varepsilon}_{g,t}, \quad (15)$$

$$\boldsymbol{\varepsilon}_{y,t} = \beta'\boldsymbol{\varepsilon}_{g,t} + \eta_t, \quad (16)$$

where $\boldsymbol{\alpha}'_g = \theta'\mathbf{B}$, $\boldsymbol{\alpha}'_y = \beta\boldsymbol{\alpha}'_g + \mathbf{v}'\mathbf{B}$, and $\boldsymbol{\gamma}' = \beta\theta' + \mathbf{v}'$. The model contains high-dimensional latent controls \mathbf{u}_t . Here, $(\boldsymbol{\alpha}_y, \boldsymbol{\alpha}_g, \beta)$ are low-dimensional coefficient vectors while $(\boldsymbol{\gamma}, \theta)$ are high-dimensional sparse vectors. Fan, Ke, and Wang (2020) and Hansen and Liao (2018) showed that the penalized regression can be successfully applied to (14) to select components in \mathbf{u}_t , which are cross-sectionally weakly correlated. They require strong conditions so that we can consistently estimate the number of factors $r = \dim(\mathbf{f}_t)$ first.

The main result of this section is to show that the factor-augmented post-selection inference is valid for any $R \geq r$. Therefore, we have addressed an important question in empirical applications, where the evidence of the number of factors is not so strong and one may use a slightly larger number of “working factors.” The theoretical intuition, again, is that the model depends on \mathbf{f}_t only through transformation invariant terms, so that $\widehat{\boldsymbol{\alpha}}'_y \widehat{\mathbf{f}}_t = \boldsymbol{\alpha}'_y \mathbf{H}^+ \mathbf{H} \mathbf{f}_t + o_p(1) = \boldsymbol{\alpha}'_y \mathbf{f}_t + o_p(1)$. In addition, \mathbf{u}_t can also be well estimated with over-identified number of factors.

Importantly, we admit the special case $r = 0$, and $R \geq 1$, leading to $\boldsymbol{\alpha}_y$ and $\boldsymbol{\alpha}_g$ both being zero in (14). That is, there are no factors, so $\mathbf{x}_t = \mathbf{u}_t$ itself is cross-sectionally weakly dependent, but nevertheless we estimate $R \geq 1$ number of factors to run post-selection inference. This setting is empirically relevant as it allows to avoid pretesting the presence of common factors for inference. The simulation in Section 5 shows that with $R \geq r$, this procedure works well even if $r = 0$; but when $r (r \geq 1)$ factors are present, directly selecting \mathbf{x}_t leads to severely biased estimations. Therefore as a practical guidance, we recommend that one should always run factor-augmented post-selection inference, with $R \geq 1$, to guard against confounding factors among the control variables.

Below we present the factor-augmented algorithm as in Hansen and Liao (2018) for estimating (11). For notational simplicity, we focus on the univariate case $\dim(\beta) = 1$.

Algorithm 3.1. Estimate β as follows.

- Step 1: Fix the working number of factors R . Estimate $\{(\mathbf{f}_t, \mathbf{u}_t) : t \leq T\}$ as in Section 2.
- Step 2: (1) Estimate coefficients: $\widehat{\boldsymbol{\alpha}}_y = (\sum_{t=1}^T \widehat{\mathbf{f}}_t \widehat{\mathbf{f}}_t')^{-1} \sum_{t=1}^T \widehat{\mathbf{f}}_t y_t$, and $\widehat{\boldsymbol{\alpha}}_g = (\sum_{t=1}^T \widehat{\mathbf{f}}_t \widehat{\mathbf{f}}_t')^{-1} \sum_{t=1}^T \widehat{\mathbf{f}}_t \mathbf{g}_t$. (2) Run penalized regression:

$$\widetilde{\boldsymbol{\gamma}} = \arg \min_{\boldsymbol{\gamma}} \frac{1}{T} \sum_{t=1}^T (y_t - \widehat{\boldsymbol{\alpha}}'_y \widehat{\mathbf{f}}_t - \boldsymbol{\gamma}'\widehat{\mathbf{u}}_t)^2 + P_\tau(\boldsymbol{\gamma}),$$

$$\widetilde{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \frac{1}{T} \sum_{t=1}^T (\mathbf{g}_t - \widehat{\boldsymbol{\alpha}}'_g \widehat{\mathbf{f}}_t - \boldsymbol{\theta}'\widehat{\mathbf{u}}_t)^2 + P_\tau(\boldsymbol{\theta}).$$

- (3) Run post-selection refitting: let $\widehat{\mathcal{J}} = \{j \leq p : \widetilde{\gamma}_j \neq 0\} \cup \{j \leq p : \widetilde{\theta}_j \neq 0\}$.

$$\widehat{\boldsymbol{\gamma}} = \arg \min_{\boldsymbol{\gamma}} \frac{1}{T} \sum_{t=1}^T (y_t - \widehat{\boldsymbol{\alpha}}'_y \widehat{\mathbf{f}}_t - \boldsymbol{\gamma}'\widehat{\mathbf{u}}_t)^2, \quad \text{such that}$$

$$\widehat{\gamma}_j = 0 \text{ if } j \notin \widehat{\mathcal{J}}.$$

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{T} \sum_{t=1}^T (\mathbf{g}_t - \hat{\alpha}_g' \mathbf{f}_t - \theta' \mathbf{u}_t)^2, \quad \text{such that}$$

$$\hat{\theta}_j = 0 \text{ if } j \notin \hat{J}.$$

Step 3: Estimate residuals: $\hat{\mathbf{e}}_{y,t} = y_t - (\hat{\alpha}_y' \mathbf{f}_t + \hat{\gamma}' \mathbf{u}_t)$, and $\hat{\mathbf{e}}_{g,t} = \mathbf{g}_t - (\hat{\alpha}_g' \mathbf{f}_t + \hat{\theta}' \mathbf{u}_t)$.

Step 4: Estimate β by residual-regression:

$$\hat{\beta} = \left(\sum_{t=1}^T \hat{\mathbf{e}}_{g,t}^2 \right)^{-1} \sum_{t=1}^T \hat{\mathbf{e}}_{g,t} \hat{\mathbf{e}}_{y,t}.$$

One can also simplify step 2 following Fan, Ke, and Wang (2020): finding $(\hat{\alpha}_y, \hat{\gamma})$ by minimizing $\frac{1}{T} \sum_{t=1}^T (y_t - \alpha_y' \mathbf{f}_t - \gamma' \mathbf{u}_t)^2 + P_{\tau}(\gamma)$ and defining $(\hat{\alpha}_g, \hat{\theta})$ similarly. Note that $\gamma \rightarrow P_{\tau}(\gamma)$ is a sparse-induced penalty function with a tuning parameter τ . In the main theorem below, we prove for the lasso $P_{\tau}(\gamma) = \tau \|\gamma\|_1$, where $\|\gamma\|_1 = \sum_{j=1}^N |\gamma_j|$. Following Bickel, Ritov, and Tsybakov (2009), set

$$\tau = C \sqrt{\frac{\sigma^2 \log N}{T}}$$

for some constant $C > 4$, where $\sigma^2 = \text{var}(\mathbf{e}_{y,t})$ for estimating γ , and $\sigma^2 = \text{var}(\mathbf{e}_{g,t})$ for estimating θ . We refer to Belloni, Chernozhukov, and Hansen (2014) for feasible tunings so that σ^2 is estimated iteratively.

3.2.2. The Main Result

We impose the following assumptions.

Assumption 3.2.

- (i) $\mathbb{E}(\mathbf{e}_{g,t} | \mathbf{u}_t, \mathbf{f}_t, \mathbf{W}) = 0$ and $\mathbb{E}(\mathbf{e}_{y,t} | \mathbf{u}_t, \mathbf{f}_t, \mathbf{W}) = 0$,
- (ii) Coefficients: there is $C > 0$, so that $\|\alpha_y\|$, $\|\alpha_g\|$, $\|\beta\|$ are all bounded by C .
- (iii) Weak dependence: There is $C > 0$, almost surely, $\max_{s \leq T} \sum_{t \leq T} |\mathbb{E}(\mathbf{e}_{y,t} \mathbf{e}_{y,s} | \mathbf{F}, \mathbf{U}, \mathbf{W})| + \max_{s \leq T} \sum_{t \leq T} |\mathbb{E}(\mathbf{e}_{g,t} \mathbf{e}_{g,s} | \mathbf{F}, \mathbf{U}, \mathbf{W})| < C$.
- (iv) Uniform bounds: $\max_{i \leq N} \|\frac{1}{T} \sum_{t=1}^T u_{it} \mathbf{v}_t\| = O_P(\sqrt{\frac{\log N}{T}})$ for all $\mathbf{v}_t \in \{\mathbf{e}_{g,t}, \mathbf{e}_{y,t}, \mathbf{f}_t\}$. In addition, $\max_{i \leq N} |\frac{1}{T} \sum_{t=1}^T (u_{it} u_{jt} - \mathbb{E} u_{it} u_{jt})| = O_P(\sqrt{\frac{\log N}{T}})$, and $\max_{i \leq N} |\frac{1}{TN} \sum_{t=1}^T \sum_{j=1}^N (u_{it} u_{jt} - \mathbb{E} u_{it} u_{jt}) w_{k,j}| = O_P(\sqrt{\frac{\log N}{TN}})$ for all $k \leq R$.

Assumption 3.2(iv) holds generally under weak time-series dependent conditions for $\{(\mathbf{v}_t, \mathbf{u}_t) : t \leq N\}$ with sub-Gaussian tails.

Suppose the high-dimensional coefficients θ and γ are strictly sparse. Let J denote the nonzero index set:

$$J = \{j \leq N : \theta_j \neq 0\} \cup \{j \leq N : \gamma_j \neq 0\}.$$

The following *sparse eigenvalue condition* is standard for the post-selection inference. Note that it is imposed on the covariance of \mathbf{u}_t rather than \mathbf{x}_t , because \mathbf{u}_t is weakly dependent.

Assumption 3.3 (Sparse eigenvalue condition). For any $\mathbf{v} \in \mathbb{R}^N \setminus \{0\}$, define

$$\phi_{\min}(m) = \inf_{\mathbf{v} \in \mathbb{R}^N : 1 \leq \|\mathbf{v}\|_0 \leq m} \mathcal{R}(\mathbf{v}) \quad \text{and}$$

$$\phi_{\max}(m) = \sup_{\mathbf{v} \in \mathbb{R}^N : 1 \leq \|\mathbf{v}\|_0 \leq m} \mathcal{R}(\mathbf{v}),$$

where $\mathcal{R}(\mathbf{v}) := \|\mathbf{v}\|^{-2} \mathbf{v}' \frac{1}{T} \sum_{t=1}^T \mathbf{u}_t \mathbf{u}_t' \mathbf{v}$. Then there is a sequence $l_T \rightarrow \infty$ and $c_1, c_2 > 0$ so that with probability approaching one,

$$c_1 < \phi_{\min}(l_T |J|_0) \leq \phi_{\max}(l_T |J|_0) < c_2.$$

Assumption 3.4.

- (i) $\frac{1}{T} \sum_{t=1}^T \mathbf{e}_{g,t}^2 \xrightarrow{P} \sigma_g^2$ for some $\sigma_g^2 > 0$.
- (ii) $\frac{1}{\sqrt{T}} \sum_{t=1}^T \eta_t \mathbf{e}_{g,t} \xrightarrow{d} \mathcal{N}(0, \sigma_{\eta g}^2)$ for some $\sigma_{\eta g}^2 > 0$. In addition, there is a consistent variance estimator $\hat{\sigma}_{\eta g}^2 \xrightarrow{P} \sigma_{\eta g}^2$.
- (iii) The rates $(N, T, |J|_0)$ satisfy

$$|J|_0^4 \log^2 N = o(T) \quad \text{and}$$

$$T |J|_0^4 = o(N^2 \min\{1, |J|_0^4 \nu_{\min}^4(\mathbf{H})\}).$$

Condition 3.4(iii) requires the “effective dimension” $N \nu_{\min}^2(\mathbf{H})$ be relatively large to accurately estimate the latent factors.

Theorem 3.2. Suppose $\hat{\mathbf{f}}_t$ contains $R \geq r \geq 0$ number of diversified weighted averages of \mathbf{x}_t . If $r \geq 1$ (there are factors in \mathbf{x}_t), Assumptions 2.1–2.4 and 3.2–3.4 hold. If $r = 0$ (there are no factors in \mathbf{x}_t), Assumption 2.2 is relaxed, and all \mathbf{f}_t involved in the above assumptions can be removed. Then as $T, N \rightarrow \infty$, for all bounded $R \geq r \geq 0$,

$$\sigma_{\eta g}^{-1} \sigma_g^2 \sqrt{T} (\hat{\beta} - \beta) \xrightarrow{d} \mathcal{N}(0, 1).$$

Fix a significant level τ , let ζ_{τ} be the $(1 - \tau/2)$ quantile of standard normal distribution. In addition, let $\hat{\sigma}_g^2 = \frac{1}{T} \sum_{t=1}^T \hat{\mathbf{e}}_{g,t}^2$. Immediately, we have the following uniform coverage.

Corollary 3.1. Suppose the assumptions of Theorem 3.2 hold. Let $\bar{R} > 0$ be a fixed upper bound for R . Then uniformly for all $0 \leq r \leq R \leq \bar{R}$,

$$\mathbb{P} \left(\beta \in [\hat{\beta} \pm \frac{1}{\sqrt{T}} \hat{\sigma}_{\eta g} \hat{\sigma}_g^{-2} \zeta_{\tau}] \right) \rightarrow 1 - \tau.$$

The novelty of the above uniformity is that the coverage is valid uniformly for all bounded r as the true number of factors, and all over-estimated R as the working number of factors. In particular, it also admits the weak-dependence $r = 0$ while $R \geq 1$ as a special case.

Remark 3.1 (Case $r = 0, R \geq 1$). We now explain the intuition of the case $\mathbf{x}_t = \mathbf{u}_t$ (no presence of confounding factors), but we nevertheless extract $R \geq 1$ “factors.” In this case $\alpha_y = \alpha_g = 0$ in the system (14). Then $\hat{\mathbf{f}}_t = \frac{1}{N} \sum_{i=1}^N \mathbf{w}_i u_{it} := \mathbf{e}_t$ degenerates to zero. Both \mathbf{u}_t and $\alpha_y' \mathbf{f}_t$ (which is zero) are still estimated well in the following sense:

$$\max_{i \leq N} \frac{1}{T} \sum_{t=1}^T (\hat{u}_{it} - u_{it})^2 = O_P \left(\frac{1}{N} + \frac{\log N}{T} \right),$$

$$\frac{1}{T} \sum_{t=1}^T (\hat{\alpha}_y' \mathbf{f}_t)^2 = O_P \left(\frac{|J|_0^2}{N} + \frac{|J|_0^2}{T} \right).$$

Remark 3.2 (Case $R = 0$). For completeness of the theorem, we define the estimator for the case $R = 0$. In this case, we do not extract any factor estimators, and simply set $\hat{\alpha}_y = \hat{\alpha}_g = 0$, and $\hat{\mathbf{u}}_t = \mathbf{x}_t$ in Algorithm 3.1. This is then the same setting as in Belloni, Chernozhukov, and Hansen (2014).

3.3. Estimating the Idiosyncratic Covariance

The estimation of the $N \times N$ idiosyncratic covariance matrix $\Sigma_u := \mathbb{E}\mathbf{u}_t\mathbf{u}_t'$ is of general interest in many applications. Examples include the efficient estimation of factor models (Bai and Li 2012), high-dimensional testing (Fan, Liao, and Yao 2015), and bootstrapping latent factors (Goncalves and Perron 2020), among many others. While this problem has been studied by Fan, Liao, and Mincheva (2013), they require that the true number of factors r has to be either known or consistently estimated, and the factors are estimated through PCA. Here we show that using the diversified factors, their conclusion holds for all fixed $R \geq r$.

A key assumption is that $\Sigma_u = (\sigma_{u,ij})$ is sparse: As in Bickel and Levina (2008), the sparsity of Σ_u is measured by the following quantity:

$$m_N = \max_{i \leq N} \sum_{j \leq N} |\sigma_{u,ij}|^q, \quad \text{for some } q \in [0, 1].$$

Given the estimated residual \hat{u}_{it} that is obtained using a working number of factors R , we estimate $\mathbb{E}u_{it}u_{jt}$ by applying a generalized thresholding: define $s_{u,ij} := \frac{1}{T} \sum_{t=1}^T \hat{u}_{it}\hat{u}_{jt}$,

$$\hat{\sigma}_{u,ij} = \begin{cases} s_{u,ij}, & \text{if } i = j, \\ h(s_{u,ij}, \tau_{ij}), & \text{if } i \neq j, \end{cases}$$

where $h(s, \tau)$ is a thresholding function with threshold value τ . Then the sparse idiosyncratic covariance estimator is defined as $\hat{\Sigma}_u = (\hat{\sigma}_{u,ij})_{N \times N}$. The threshold value τ_{ij} is chosen as

$$\tau_{ij} = C\sqrt{s_{u,ii}s_{u,jj}\omega_{NT}}, \quad \omega_{NT} := \sqrt{\frac{\log N}{T}} + \frac{1}{\sqrt{N}}$$

for some large constant $C > 0$, which applies a constant thresholding to correlations.

In general, the thresholding function should satisfy:

- (i) $h(s, \tau) = 0$ if $|s| < \tau$,
- (ii) $|h(s, \tau) - s| \leq \tau$.
- (iii) there are constants $a > 0$ and $b > 1$ such that $|h(s, \tau) - s| \leq a\tau^2$ if $|s| > b\tau$.

Note that condition (iii) requires that the thresholding bias should be of higher order. It is not necessary for consistent estimations, but we recommend using nearly unbiased thresholding (Antoniadis and Fan 2001) for inference applications. One such example is known as SCAD. As noted in Fan, Liao, and Yao (2015), the unbiased thresholding is required to avoid size distortions in a large class of high-dimensional testing problems involving a “plug-in” estimator of Σ_u . In particular, this rules out the popular *soft-thresholding* function, which does not satisfy (iii) due to its first-order shrinkage bias.

Theorem 3.3. Let $\hat{\mathbf{u}}_t$ be constructed using $R \geq r$ number of diversified weighted averages of \mathbf{x}_t . Suppose that Assumptions 2.1–2.4 hold and that $\log N = o(T)$. In addition, either $v_{\min}^2(\mathbf{H}) \gg \frac{1}{\sqrt{N}}$ or $v_{\min}^2(\mathbf{H}) \gg \frac{1}{N}\sqrt{\frac{T}{\log N}}$. Then as $N, T \rightarrow \infty$, for any $R \geq r \geq 0$,

(i)

$$\max_{i \leq N} \frac{1}{T} \sum_{t=1}^T (\hat{\mathbf{b}}_i' \hat{\mathbf{f}}_t - \mathbf{b}_i' \mathbf{f}_t)^2 = O_P(\omega_{NT}).$$

(ii) For a sufficiently large constant $C > 0$ in the threshold τ_{ij} ,

$$\|\hat{\Sigma}_u - \Sigma_u\| = O_P(\omega_{NT}^{1-q} m_N).$$

(iii) If in addition, $\lambda_{\min}(\Sigma_u) > c_0$ for some $c_0 > 0$ and $\omega_{NT}^{1-q} m_N = o(1)$, then

$$\|\hat{\Sigma}_u^{-1} - \Sigma_u^{-1}\| = O_P(\omega_{NT}^{1-q} m_N).$$

3.4. Testing Specification of Factors

In practical applications, many “observed factors” \mathbf{g}_t have been proposed to approximate the true latent factors. For example, in asset pricing, popular choices of \mathbf{g}_t are proposed and discussed in seminal works by Fama and French (1992) and Carhart (1997), which are known as the Fama–French factors and Carhart four factor models.

We test the (linear) specification of a given set of empirical factors \mathbf{g}_t . That is, we test:

$$H_0: \text{there is a } r \times r \text{ invertible matrix } \boldsymbol{\theta} \text{ so that } \mathbf{g}_t = \boldsymbol{\theta} \mathbf{f}_t, \\ \forall t \leq T.$$

Under the null hypothesis, \mathbf{g}_t and \mathbf{f}_t are linear functions of each other. We propose a simple statistic:

$$\|\mathbf{P}_G - \mathbf{P}_{\hat{\mathbf{F}}}\|_F^2,$$

where $\mathbf{G} = (\mathbf{g}_1, \dots, \mathbf{g}_T)'$ and recall that $\mathbf{P}(\cdot)$ denotes the projection matrix. Here we still use the diversified factor estimator $\hat{\mathbf{F}}$. The test statistic measures the distance between (linear) spaces, respectively, spanned by \mathbf{g}_t and $\hat{\mathbf{f}}_t$. To derive the asymptotic null distribution, we naturally set the working number of factors $R = \dim(\mathbf{g}_t)$, which is known and equals $\dim(\mathbf{f}_t) = r$ under the null. Then $\|\mathbf{P}_{\hat{\mathbf{F}}} - \mathbf{P}_{\mathbf{F}}\|_F = o_P(1)$, followed from Theorem 2.1.

3.4.1. Asymptotic Null Distribution

With the diversified factor estimators, the null distribution of the statistic is very easy to derive, and satisfies:

$$\frac{N\sqrt{T}(\|\mathbf{P}_G - \mathbf{P}_{\hat{\mathbf{F}}}\|_F^2 - \text{MEAN})}{\sigma} \xrightarrow{d} \mathcal{N}(0, 1),$$

where for $\mathbf{A} = 2\mathbf{H}'^{-1}(\frac{1}{T}\mathbf{F}'\mathbf{F})^{-1}\mathbf{H}^{-1}$,

$$\text{MEAN} = \frac{1}{N^2} \text{tr} \mathbf{A} \mathbf{W}' \mathbb{E}(\mathbf{u}_t \mathbf{u}_t' | \mathbf{F}) \mathbf{W},$$

$$\sigma^2 = \text{Var} \left(\frac{1}{N} \text{tr} \mathbf{A} \mathbf{W}' \mathbf{u}_t \mathbf{u}_t' \mathbf{W} | \mathbf{F}, \mathbf{W} \right) > 0.$$

Here, we assume $\sigma^2 > 0$ to be bounded away from zero. To avoid nonparametrically estimating high-dimensional covariances, we shall assume the conditional covariances in both bias and variance are independent of \mathbf{F} almost surely. Nevertheless, the bias depends on a high-dimensional matrix $\Sigma_u = \mathbb{E}(\mathbf{u}_t \mathbf{u}_t')$. We employ the sparse covariance $\widehat{\Sigma}_u$ as defined in Section 3.3 and replace the bias by

$$\widehat{\text{MEAN}} := \frac{1}{N^2} \text{tr} \widehat{\mathbf{A}} \mathbf{W}' \widehat{\Sigma}_u \mathbf{W} \quad \text{with} \quad \widehat{\mathbf{A}} := 2 \left(\frac{1}{T} \widehat{\mathbf{F}} \widehat{\mathbf{F}}' \right)^{-1}.$$

Further suppose σ can be consistently estimated by some $\widehat{\sigma}$, then together, we have the feasible standardized statistic:

$$\frac{N\sqrt{T}(\|\mathbf{P}_{\widehat{\mathbf{F}}} - \mathbf{P}_{\mathbf{G}}\|_F^2 - \widehat{\text{MEAN}})}{\widehat{\sigma}}. \quad (17)$$

The problem, however, is not as straightforward as it looks by far. The use of $\widehat{\text{MEAN}}$ and $\widehat{\sigma}$ both come with issues, as we now explain.

3.4.1.1. The Issue of $\widehat{\text{MEAN}}$. When deriving the asymptotic null distribution, we need to address the effect of $\widehat{\Sigma}_u - \Sigma_u$, which is to show

$$\begin{aligned} & \frac{N\sqrt{T}(\widehat{\text{MEAN}} - \text{MEAN})}{\sigma} \\ & \approx \frac{N\sqrt{T}}{\sigma} \frac{1}{N^2} \text{tr} \mathbf{A} \mathbf{W}' (\widehat{\Sigma}_u - \Sigma_u) \mathbf{W} \xrightarrow{P} 0. \end{aligned} \quad (18)$$

But simply applying the rate of convergence of $\|\widehat{\Sigma}_u - \Sigma_u\|$ in Theorem 3.3 fails to show the above convergence, even though the rate is minimax optimal³. Similar phenomena also arise in Fan, Liao, and Yao (2015) and Bai and Liao (2017), where a plug-in estimator for Σ_u is used for inferences. Proving (18) requires a dedicated technical argument to address the accumulation of high-dimensional estimation errors. It requires a strengthened condition on the weak cross-sectional dependence, in Assumption 3.8 below.

3.4.1.2. The Issue of $\widehat{\sigma}$. It is difficult to estimate σ through residuals $\widehat{\mathbf{u}}_t$ since $\mathbf{W}' \widehat{\mathbf{u}}_t = 0$ almost surely. In fact, estimated \mathbf{u}_t constructed based on any factor estimator would lead to *inconsistent* estimator for σ^2 . Therefore, we propose to estimate σ^2 by parametric bootstrap. Observe that $\frac{1}{\sqrt{N}} \mathbf{W}' \mathbf{u}_t$ is asymptotically normal, whose variance is given by $\mathbf{V} = \frac{1}{N} \mathbf{W}' \Sigma_u \mathbf{W}$. Hence, σ^2 should be approximately equal to

$$f(\mathbf{A}, \mathbf{V}) := \text{var} \left(\frac{1}{N} \text{tr} \mathbf{A} \mathbf{W}' \mathbf{Z}_t \mathbf{Z}_t' \mathbf{W} \right), \quad (19)$$

where \mathbf{Z}_t is distributed as $\mathcal{N}(0, \mathbf{V})$. Therefore, we estimate σ^2 by

$$\widehat{\sigma}^2 = f(\widehat{\mathbf{A}}, \widehat{\mathbf{V}}), \quad \text{with} \quad \widehat{\mathbf{V}} = \frac{1}{N} \mathbf{W}' \widehat{\Sigma}_u \mathbf{W},$$

which can be calculated by simulating from $\mathcal{N}(0, \widehat{\mathbf{V}})$.

Above all, despite of the simple construction of $\widehat{\mathbf{F}}$, the technical problem is still challenging. Therefore, this subsection calls for relatively stronger conditions, as we now impose.

Assumption 3.5.

- (i) $\{\mathbf{u}_t : t \leq T\}$ are stationary and conditionally serially independent, given \mathbf{F} and \mathbf{G} .
- (ii) There is $C > 0$, $\mathbb{E}[\|\frac{1}{\sqrt{N}} \mathbf{W}' \mathbf{u}_t\|^4 | \mathbf{W}] < C$.
- (iii) $\nu_{\min}(\mathbf{H}) > c$ for some $c > 0$.

The next assumption ensures that σ^2 can be estimated by simulating from the Gaussian distribution.

Assumption 3.6.

- (i) There is $c > 0$ so that $\sigma^2 > c$.
- (ii) As $N \rightarrow \infty$, $|\sigma^2 - f(\mathbf{A}, \mathbf{V})| \rightarrow 0$ almost surely in \mathbf{F} , where $f(\mathbf{A}, \mathbf{V})$ is given in (19).

Next, we shall require Σ_u be strictly sparse, in the sense that the “small” off-diagonal entries are exactly zero. In this case, we use the following measurement for the total sparsity:

$$D_N := \sum_{i,j \leq N} 1\{\mathbb{E} u_{it} u_{jt} \neq 0\}.$$

Recall that $\omega_{NT} := \sqrt{\frac{\log N}{T}} + \frac{1}{\sqrt{N}}$. We assume:

Assumption 3.7 (Strict sparsity).

- (i) $(\frac{\omega_{NT}^2 \sqrt{T}}{N}) D_N \rightarrow 0$.
- (ii) $\min\{|\mathbb{E} u_{it} u_{jt}| : \mathbb{E} u_{it} u_{jt} \neq 0\} \gg \omega_{NT}$.

For block-diagonal matrices with finite block sizes, $D_N = O(N)$; for banded matrices with band size l_N , $D_N = O(l_N N)$. In general, suppose $D_N = l_N N$ with some slowly growing $l_N \rightarrow \infty$. Then condition (i) reduces to requiring $l_N^2 \log N \ll l_N \sqrt{T} \ll N$. This requires an upper bound for l_N ; in addition, the lower bound for N arises from the requirement of estimating factors. Condition (ii) requires that the nonzero entries are well-separated from the statistical errors.

Assumption 3.8. Write $\sigma_{u,ij} := \mathbb{E} u_{it} u_{jt}$. There is $C > 0$ so that

$$\frac{1}{N} \sum_{(m,n): \sigma_{u,mn} \neq 0} \sum_{(i,j): \sigma_{u,ij} \neq 0} |\text{cov}(u_{it} u_{jt}, u_{mt} u_{nt})| < C.$$

The above assumption is the key condition to argue for (18). It requires further conditions on the weak cross-sectional dependence, in addition to the sparsity. Fan, Liao, and Yao (2015) proved that if u_{it} is Gaussian, then a sufficient condition for Assumption 3.8 is as follows:

$$D_N = O(N), \text{ and } \max_{i \leq N} \sum_{j \leq N} 1\{\mathbb{E} u_{it} u_{jt} \neq 0\} = O(1),$$

which is the case for block diagonal matrices with finite members in each block and banded matrices with $l_N = O(1)$.

Theorem 3.4. Suppose $R = \dim(\mathbf{g}_t)$, and Assumptions 2.1–2.4 and 3.5–3.8 hold. As $N, T \rightarrow \infty$, under H_0 ,

$$\frac{N\sqrt{T}(\|\mathbf{P}_{\widehat{\mathbf{F}}} - \mathbf{P}_{\mathbf{G}}\|_F^2 - \widehat{\text{MEAN}})}{\widehat{\sigma}} \xrightarrow{d} \mathcal{N}(0, 1).$$

³A simple calculation would only yield $\frac{N\sqrt{T}}{\sigma} \frac{1}{N^2} \|\mathbf{A} \mathbf{W}'\| \|\widehat{\Sigma}_u - \Sigma_u\| \|\mathbf{W}\| \leq O_p(1)$ but not necessarily $o_p(1)$.

3.5. Factor-Adjusted False Discovery Control for Multiple Testing

Controlling the false discovery rate (FDR) in large-scale hypothesis testing based on strongly correlated testing series has been an important problem. Suppose the data are generated from:

$$\mathbf{x}_t = \boldsymbol{\alpha} + \mathbf{B}\mathbf{f}_t + \mathbf{u}_t,$$

where $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_N)'$ is the mean vector. This model allows strong cross-sectional dependences among \mathbf{x}_t . We are interested in testing N number of hypotheses:

$$H_0^i : \alpha_i = 0, \quad i = 1, \dots, N.$$

The FDR control aims to develop test statistics Z_i and threshold values so that the overall FDR is controlled at certain value. A crucial requirement is that these test statistics should be weakly dependent. However, for $\bar{\mathbf{f}} = \frac{1}{T} \sum_t \mathbf{f}_t$ and $\bar{\mathbf{u}} = \frac{1}{T} \sum_t \mathbf{u}_t$, we have $\bar{\mathbf{x}} = \frac{1}{T} \sum_t \mathbf{x}_t = \boldsymbol{\alpha} + \mathbf{B}\bar{\mathbf{f}} + \bar{\mathbf{u}}$, so the presence of $\mathbf{B}\mathbf{f}_t$ makes the mean vector be cross-sectionally strongly dependent, failing usual FDR procedures based on the simple sample average. This is the well-known confounding factor problem. While several methods have been proposed to remove the effect of confounding factors (Wang et al. 2017; Fan et al. 2019), again, it has been assumed that the number of factors should be consistently estimable.

The diversified projection can be applied directly as a simple implementation for the FDR control, valid for all $R \geq r$. Let the diversified projection be $\hat{\mathbf{f}}_t = \frac{1}{N} \mathbf{W}' \mathbf{x}_t$, and let $\hat{\mathbf{b}}_i$ be the OLS estimator for the slope vector by regressing x_{it} on $\hat{\mathbf{f}}_t$ with intercept. Then we can define the factor-adjusted regularized multiple test (Fan et al. 2019) statistics $Z_i = \hat{\alpha}_i / \text{SE}(\hat{\alpha}_i)$ where

$$\hat{\alpha}_i = \bar{x}_i - \hat{\mathbf{b}}_i' \bar{\mathbf{f}}, \quad \hat{\mathbf{f}} = \frac{1}{T} \sum_{t=1}^T \hat{\mathbf{f}}_t,$$

and $\text{SE}(\hat{\alpha}_i)$ is the associated standard error. Our theories imply the following expansion, uniformly for $i = 1, \dots, N$ and all $R \geq r$,

$$\hat{\alpha}_i - \alpha_i = \frac{1}{T} \sum_{t=1}^T \mathbf{g}_t' \mathbf{u}_{it} + o_p(T^{-1/2}),$$

where $\mathbf{g}_t = 1 - \bar{\mathbf{f}}' \mathbf{S}_f^{-1} (\mathbf{f}_t - \bar{\mathbf{f}})$, and $\mathbf{S}_f = \frac{1}{T} \sum_t (\mathbf{f}_t - \bar{\mathbf{f}})(\mathbf{f}_t - \bar{\mathbf{f}})'$. This gives rise to the desired expansion so that Z_i are weakly dependent. Therefore, we can apply standard procedures to Z_i for the false discovery control.

4. Choices of Diversified Weights

We discuss some specific examples to choose $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_R) = (w_{k,i} : k \leq R, i \leq N)$, the weight matrix.

4.1. Loading Characteristics

Factor loadings are often driven by observed characteristics. For example, in genetic studies, single-nucleotide polymorphism (SNP) data are often collected with the gene expression data on the same group of subjects. The SNPs drive underlying

structure in the gene expressions, clinical and demographics data, through affecting their loadings on the biological factors. In asset pricing studies, it has been well documented that factor loadings are driven by firm specific characteristics, which are independent of the model noise, but have strong explanatory powers on the loadings.

Motivated by the presence of characteristics, “characteristic based” factor models have been extensively studied in the literature (e.g., Connor, Matthias, and Linton 2012; Gagliardini, Ossola, and Scaillet 2016; Li et al. 2016). The general form of this model assumes the loadings have the following decomposition (Fan, Liao, and Wang 2016):

$$\mathbf{b}_i = \mathbf{g}(\mathbf{z}_i) + \boldsymbol{\gamma}_i, \quad \mathbb{E}(\boldsymbol{\gamma}_i | \mathbf{z}_i) = 0, \quad i \leq N, \quad (20)$$

where \mathbf{z}_i is a vector of characteristics that are observed on each subject and $\mathbf{g}(\cdot)$ is a nonparametric mean function. It is assumed that $\{\mathbf{z}_i : i \leq N\}$ is independent of \mathbf{u}_t and that $\mathbf{g}(\mathbf{z}_i)$ is not degenerate so that \mathbf{z}_i has explanatory power. In addition, $\boldsymbol{\gamma}_i$ is the remaining loading components, after conditioning on \mathbf{z}_i . The decomposition of \mathbf{b}_i in (20) is motivated from the asset pricing literature, where factor “betas” are known to be partially explained by individual-specific observables \mathbf{z}_i , which represent a set of time-invariant characteristics such as individual stocks’ size, momentum, values. When \mathbf{z}_i is available, we can employ them as a natural choice of the weights for the diversified factors. Fix an R -component of sieve basis functions: $(\phi_1(\cdot), \dots, \phi_R(\cdot))$ such as the Fourier basis or B splines. Then define

$$\mathbf{W} := (w_{i,k})_{N \times R}, \quad \text{where } w_{i,k} = \phi_k(\mathbf{z}_i).$$

The diversified projection using the so-constructed \mathbf{W} is related to the “projected PCA” of Fan, Liao, and Wang (2016), but the latter is more complicated and requires stronger conditions than the diversified projection, because it is still PCA based.

4.2. Moving Window Estimations

This method is useful when \mathbf{u}_t is serially independent, and related ideas have been used recently by Barigozzi and Cho (2018). Consider out-of-sample forecasts using moving windows. Suppose \mathbf{x}_t is observed for $T + T_0$ periods in total, but to pertain the stationarity assumption, we only use the most recent T observations to learn the latent factors, where T may be potentially small. Divide the sample into two periods:

periods (I) of learning weights: $\mathbf{x}_t = \mathbf{B}_1 \mathbf{f}_t + \mathbf{u}_t$,

$$t = 1, \dots, T_0,$$

periods (II) of interest: $\mathbf{x}_t = \mathbf{B} \mathbf{f}_t + \mathbf{u}_t$,

$$t = T_0 + 1, \dots, T_0 + T.$$

While \mathbf{B}_1 and \mathbf{B} can be different (e.g., presence of structural breaks), they are assumed to be closely related between two sampling periods. As such, we can learn about the diversified weights from periods (I) to estimate the latent factors for the periods of estimation interests (II). Specifically, apply PCA on periods (I) to extract R number of factor loadings:

$\widehat{\mathbf{B}}_1 = (\widehat{b}_{i,k})_{N \times R}$. Now for a predetermined constant $\epsilon > 0$, define $\mathbf{W} = (w_{i,k})_{N \times R}$ where

$$w_{i,k} = \frac{\widehat{b}_{i,k}}{\max\{1, \epsilon \max_{i \leq N} |\widehat{b}_{i,k}|\}}, \quad k \leq R, \quad i \leq N.$$

Barigozzi and Cho (2018) suggested a specific choice for ϵ that work well in their simulation studies. As discussed by these authors, the trimming constant ϵ ensures that with a large probability most of the corresponding $\widehat{b}_{i,k}$ are “preserved” by $w_{i,k}$. On the other hand, if a few elements of $\widehat{b}_{i,k}$ are spuriously large in finite sample, the trimming shrinks the corresponding $\widehat{b}_{i,k}$ downward to $1/\epsilon$.

In addition, if \mathbf{u}_t is serially independent, then \mathbf{W} is also independent of \mathbf{u}_t for $t = m+1, \dots, m+T$. As such, the conditions on the diversified weights are satisfied. It is straightforward to extending this idea to multi-periods rolling window forecasts, where weights are sequentially updated for rolling windows.

The aforementioned method uses the idea that sample splitting creates serial independences. In the presence of mixing-type serial dependences, Barigozzi and Cho (2018) proposed to split the data into blocks and estimate factor loadings using subsamples omitting the current block as well as its immediate neighbors. Their method can be also applied in the current context to create the weighting matrix.

4.3. Initial Transformation

A related idea is to use transformations of the initial observation \mathbf{x}_t for $t = 0$. Suppose $(\mathbf{f}_0, \mathbf{u}_0)$ is independent of $\{\mathbf{u}_t : t \geq 1\}$, and let $\{\phi_k : k = 1, \dots, R\}$ be a set of sieve transformations. Then we can apply $w_{i,k} = \phi_k(x_{i,0})$. These weights are correlated with \mathbf{B} through $\mathbf{x}_0 = \mathbf{B}\mathbf{f}_0 + \mathbf{u}_0$ so that the rank condition is satisfied. The initial transformation method only requires $\{\mathbf{u}_t\}$ be independent of its initial value. The similar idea has been used recently by Juodis and Sarafidis (2020).

4.4. Hadamard Projection

We can set deterministic weights as in the statistical experimental designs:

$$\mathbf{W} = \begin{pmatrix} 1 & 1 & 1 & 1 & \dots \\ 1 & -1 & 1 & 1 & \dots \\ 1 & 1 & -1 & 1 & \dots \\ 1 & -1 & -1 & -1 & \dots \\ 1 & 1 & 1 & -1 & \dots \\ 1 & -1 & 1 & -1 & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix}.$$

So for each $2 \leq k \leq R$, the k th column of \mathbf{W} equals $(1'_{k-1}, -1'_{k-1}, 1'_{k-1}, -1'_{k-1}, \dots)$, where 1_m denotes the m -dimensional vector of ones. Closely related types of matrices are known as the Walsh-Hadamard matrices, formed by rearranging the columns so that the number of sign changes in a column is in an increasing order, and the columns are orthogonal. Therefore, we can also set \mathbf{W} as the $N \times R$ upper-left corner submatrix of a Hadamard matrix of dimension 2^K with $K = \lceil \log_2 N \rceil$, where $\lceil \cdot \rceil$ denotes the ceiling function.

5. Monte Carlo Experiments

In this section, we illustrate the finite sample properties of the forecasting and inference methods based on diversified factors, and use four types of weight matrices:

- (i) Hadamard weight: $\mathbf{w}_1 = \mathbf{1}$ and $\mathbf{w}_k = (\mathbf{1}'_{k-1}, -\mathbf{1}'_{k-1}, \mathbf{1}'_{k-1}, -\mathbf{1}'_{k-1}, \dots)$ for $2 \leq k \leq R$, where $\mathbf{1}_{k-1}$ is a vector of ones of length $k-1$.
- (ii) Loading characteristics: loadings depend on some characteristics \mathbf{z}_i , and we apply the polynomial transformations so that the i th row of \mathbf{W} is $(g_1(\mathbf{z}_i), g_2(\mathbf{z}_i), \dots, g_R(\mathbf{z}_i))$ for $i \leq N$. In our numerical work, we take one characteristic and set $g_j(\mathbf{z}_i) = \mathbf{z}_i^j$.
- (iii) Rolling windows: when conducting simulations for out-of-sample forecasts, we use the trimmed PCA as described in Section 4.2.
- (iv) Initial transformations: we use the initial transformation so that the i th row of \mathbf{W} is $(x_{i,0}, x_{i,0}^2, \dots, x_{i,0}^R)$ for $i \leq N$.

We generate the data from the following model motivated from Section 4.1:

$$\mathbf{x}_t = \mathbf{B}\mathbf{f}_t + \mathbf{u}_t, \quad \mathbf{B} = (b_{i,k}) * N^{-(1-\alpha)/2}, \quad \text{with} \\ b_{i,k} = (\mathbf{z}_i^k + 0.5\gamma_{i,k}).$$

We set $\mathbf{z}_i = \sin(h_i)$ where both h_i and $\gamma_{i,k}$ are independent scalar standard normal variables. Here, we use the polynomial transformation \mathbf{z}_i^k to represent the effect of characteristics. In addition, the $\gamma_{i,k}$ -component captures the unobservable beta components that are not explainable by the characteristics. With the identification condition $\mathbb{E}(\gamma_{i,k}|\mathbf{z}_i) = 0$, both components in $b_{i,k}$ can be consistently estimated. See more motivations of this model in Fan, Liao, and Wang (2016) and Kim, Korajczyk, and Neuhierl (2018). The multiplier $N^{-(1-\alpha)/2}$ measures the strength of the factors, whereas the spiked eigenvalue of the sample covariance grows at rate N^α . Hence, larger α indicates stronger factors.

The factors are multivariate standard normal. To generate the idiosyncratic term, we set the $N \times T$ matrix $\mathbf{U} = \Sigma_N^{1/2} \bar{\mathbf{U}} \Sigma_T^{1/2}$; here $\bar{\mathbf{U}}$ is an $N \times T$ matrix, whose entries independent standard normal. The $N \times N$ matrix Σ_N and the $T \times T$ matrix Σ_T , respectively, govern the cross-sectional and serial correlations of u_{it} . We set $\Sigma_T = (\rho_T^{|t-s|})_{st}$, and use a sparse cross-sectional covariance:

$$\Sigma_N = \text{diag}\{\underbrace{\mathbf{A}, \dots, \mathbf{A}}_{n \text{ of them}}, \mathbf{I}\}, \quad \mathbf{A} = (\rho_N^{|i-j|}), \quad (21)$$

where \mathbf{A} is a small four-dimensional block matrix and \mathbf{I} is $(N - 4n) \times (N - 4n)$ identity matrix so that Σ_N has a block-diagonal structure. We fix $n = 3$ and $\rho_N = 0.7$. The numerical performances are studied in the following subsection with various choice of ρ_T to test about the sensitivity against serial correlations.

5.1. Covariance Estimation

We first study the performance of estimating Σ_u . To do so, we set $r = 1$ and, respectively, calculate $\widehat{\Sigma}_u$ using $R = r, \dots, r + 3$. As estimating Σ_u is particularly important in asset pricing

models, we use the loading characteristic weights $w_{i,k} = \mathbf{z}_i^k$, $k = 1, \dots, R$, as the characteristic \mathbf{z}_i is often directly observable along with the return data.

For comparison purposes, we also estimate Σ_u using two benchmark estimators:

- (i) The PC-estimator for factors with $R = r$ (the POET method by Fan, Liao, and Mincheva (2013)). So the PC-estimator in this simulation assumes the true number of factors $r = 1$ to be known;
- (ii) The known-factor method. We use the true factors, and estimate loadings and u_{it} by OLS, followed by SCAD-thresholding.

We set two serial dependence scenarios: $\rho_T = 0.1$ (weak serial dependence) and $\rho_T = 0.7$ (strong serial dependence), as well as two factor-strength scenarios: $\alpha = 1$ and $\alpha = 0.5$.

Figure 1 plots $\|\widehat{\Sigma}_u - \Sigma_u\|$ and $\|\widehat{\Sigma}_u^{-1} - \Sigma_u^{-1}\|$, averaged over 100 replications, as $N = T$ grows. While all estimators perform similarly, the POET-estimator is not always better than the diversifying projection (DP). For estimating Σ_u , both the DP with $R = r$ and the known factor method are overall better than the POET estimator, followed by DP with other choices of R . This comparison is reasonable, reflecting the robustness of DP to the serial conditions and strength of factors. Perhaps what is surprising is the comparison for estimating the inverse covariance. In all four scenarios of the factor strength and serial correlations, the DP with $R = r$ performs the worst among the six estimators, and DP with over estimated R is in general better than both the known factor method and the POET. Our interpretation of this is that we set relatively strong cross-sectional correlations in the data-generating process, making Σ_u^{-1} more unstable. The use of more diversified weights provides extra information to help stabilizing the inverse covariance estimator.

5.2. Out-of-Sample Forecast

We assess the performance of the proposed factor estimators on out-of-sample forecasts. Consider the following forecast model

$$y_{t+1} = \beta_0 + \beta y_t + \alpha' \mathbf{f}_t + \varepsilon_{t+1},$$

where we set $r = \dim(\mathbf{f}_t) = 2$, $(\beta_0, \beta) = (1.5, 0.5)$, and $\alpha = (1, 1)'$. In addition, ε_t are independent standard normal. The data generating process for $\mathbf{x}_t = \mathbf{B}\mathbf{f}_t + \mathbf{u}_t$ is the same as before, in the presence of both serial and cross-sectional correlations. We conduct one-step-ahead out-of-sample forecast m times using a moving window of size T . Here, T is also the sample size for estimations. We simulate $m + T$ observations in total. For each $t = 0, \dots, m-1$, we use the data $\{(\mathbf{x}_{t+1}, y_{t+1}), \dots, (\mathbf{x}_{t+T}, y_{t+T})\}$ to conduct one-step-ahead forecast of y_{t+T+1} . Specifically, we estimate the factors using $\{\mathbf{x}_{t+1}, \dots, \mathbf{x}_{t+T}\}$, and obtain $\{\widehat{\mathbf{f}}_{t+1}, \dots, \widehat{\mathbf{f}}_{t+T}\}$. The coefficients in the forecasting regression is then estimated by the OLS, denoted by $(\widehat{\beta}_{0,t+T}, \widehat{\beta}_{t+T}, \widehat{\alpha}_{t+T})'$. We then forecast y_{t+T+1} by

$$\widehat{y}_{t+T+1|t+T} = \widehat{\beta}_{0,t+T} + \widehat{\beta}_{t+T} y_{t+T} + \widehat{\alpha}_{t+T}' \widehat{\mathbf{f}}_{t+T}.$$

Such a procedure continues for $t = 0, \dots, m-1$.

We compute the diversified factor estimators using the two types of weights, with $R = r, r+1, r+3$ as the working number

of factors. As for the moving windows weight, we assume there is a historical time series $\mathbf{x}_t = \mathbf{B}_1 \mathbf{f}_t + \mathbf{u}_t$, for $t = -T, \dots, 0$, and the loadings \mathbf{B}_1 is correlated with \mathbf{B} in the sense that $\mathbf{B}_1 = 0.8\mathbf{B} + 0.5\mathbf{Z}$, where \mathbf{Z} is multivariate standard normal. We then apply the moving window method to create \mathbf{W} as outlined in Section 4.2. Though the theory for the moving window weights requires serial correlation $\rho_T = 0$, we nevertheless set $\rho_T = 0, 0.5$, and 0.9 to examine the performance under serially correlated series.

The benchmark method is the PC-estimator, which uses the true number of factors. In addition, we also consider two well-known methods that specifically estimate factor dynamics:

- (i) GDF: the generalized dynamic factor model of Forni et al. (2005). The selection criterion of Hallin and Liška (2007) recommended using, on average, three dynamic factors, so we use $R = 3, 4$ numbers of dynamic factors.
- (ii) KF: the two-step Kalman filtering of Doz, Giannone, and Reichlin (2011). In the first step factors are preliminarily estimated and fit a VAR model; in the second step, Kalman smoother is applied to calculate the projection onto the observations. For this approach, we use $R = 2$, the true number of factors.

For each method M , we calculate the mean squared out-of-sample forecasting error:

$$\text{MSE}(M) = \frac{1}{m} \sum_{t=0}^{m-1} (y_{t+T+1} - \widehat{y}_{t+T+1|t+T})^2,$$

and report the relative MSE to the PC method: $\text{MSE}(M)/\text{MSE}(\text{PC})$. It is worthwhile to emphasize that this study does not aim to beat the PC-method. In fact, the PC-estimator yields the optimal rank r -estimation of the low-rank structure, in the sense that the estimated low-rank component $\mathbf{B}\mathbf{F}'$ satisfies: $\widehat{\mathbf{B}}_{pc} \widehat{\mathbf{F}}_{pc}' = \arg \min_{\text{rank}(\mathbf{A})=r} \|\mathbf{X} - \mathbf{A}\|_F^2$. So when the number of factors r is correctly specified and the time series dependence is not strong, the PC-estimator enjoys some optimal property. Nevertheless, we use PC as the benchmark as it is the most commonly used in this literature. We aim to see how well the proposed DP method performs relative to the benchmark.

The results are reported in Table 1 for $m = 50$, and is computed based on one set of simulation replications. We see that the DP with various R and generalized DF are in most scenarios similar to the PC-estimator, and DP outperforms under the strong serial correlations. In all cases, Kalman filtering is comparable with PC, including the case of strong serial correlations.

5.3. Post-Selection Inference

We now study the inference for the effect of \mathbf{g}_t in the following factor-augmented model

$$y_t = \beta \mathbf{g}_t + \mathbf{v}' \mathbf{x}_t + \eta_t,$$

$$\mathbf{g}_t = \boldsymbol{\theta}' \mathbf{x}_t + \varepsilon_{g,t},$$

$$\mathbf{x}_t = \mathbf{B} \mathbf{f}_t + \mathbf{u}_t,$$

where both \mathbf{v} and $\boldsymbol{\theta}$ are set to high-dimensional sparse vectors. The goal is to make inference about β , using the factor-augmented post-selection inference. We generate

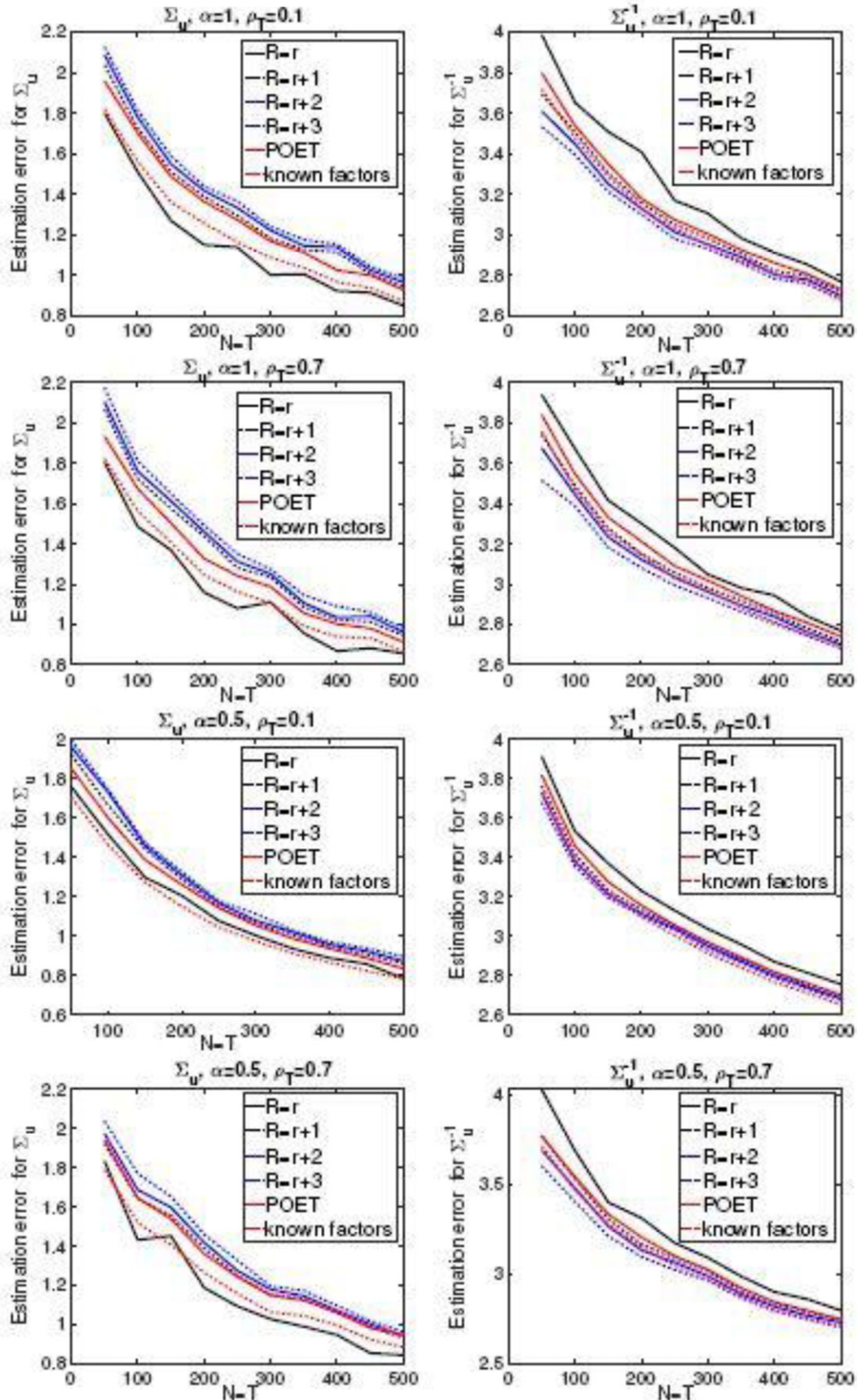


Figure 1. The estimation errors in operator-norm $\|\hat{\Sigma}_U - \Sigma_U\|$ (left) and $\|\hat{\Sigma}_U^{-1} - \Sigma_U^{-1}\|$ (right) as the dimension increases, averaged over 100 replications. We set $N = T$. Here $R = r, \dots, r+3$ correspond to the diversified factor estimators using R number of working factors. Characteristic weights are used. Here, α measures the factor strength and ρ_T is the serial correlation.

Table 1. Out-of-sample MSE(M)/MSE(PC) for three types of estimators.

ρ_T	N	T	Characteristic weights			Rolling window weights, R			GDF		KF
			r	$r + 1$	$r + 3$	r	$r + 1$	$r + 3$	3	4	r
$\alpha = 1$											
0	100	50	1.141	1.090	1.109	0.968	1.001	1.010	0.991	1.016	1.007
		100	0.998	0.980	1.035	0.979	1.039	1.046	1.008	1.009	1.002
0.5		50	0.996	1.008	0.965	0.993	1.018	1.055	1.000	0.996	0.986
		100	0.885	0.886	0.917	0.937	0.922	0.939	0.995	0.997	1.005
0.9		50	0.602	0.621	0.637	0.608	0.620	0.680	0.763	0.772	1.023
		100	0.434	0.458	0.482	0.422	0.419	0.450	0.863	0.578	0.985
$\alpha = 0.2$											
0		50	0.876	0.913	0.987	1.072	1.059	1.071	0.991	0.985	1.003
		100	0.931	0.906	0.966	1.065	1.114	1.156	0.996	1.012	0.992
0.5		50	0.891	0.897	1.044	1.059	1.082	1.149	1.002	0.981	0.958
		100	0.972	0.963	0.970	0.868	0.793	0.817	0.968	0.981	1.007
0.9		50	0.478	0.513	0.647	0.713	0.731	0.688	0.953	0.745	0.966
		100	0.762	0.765	0.767	0.788	0.806	0.849	0.927	0.851	0.951

NOTE: Reported are the out-of-sample relative MSEs. The benchmark PC-estimator uses the true number of factors. The dimension $N = 100$ is fixed. The diversified projection uses R estimated factors with two types of weights: characteristic weights and rolling window weights. In addition, the columns of GDF estimates factors from the generalized dynamic factor model of Forni et al. (2005), with R number of dynamic factors. The Matlab codes for implementing Forni et al. (2005) and Hallin and Liška (2007) are downloaded from Matteo Barigozzi's website (www.barigozzi.eu/codes.html). The column of KF refers to the Kalman filtering developed by Doz, Giannone, and Reichlin (2011), which uses the true r number of factors. Both GDF and KF specifically estimate dynamic factors.

$\mathbf{u}_t \sim \mathcal{N}(0, \Sigma_u)$, $(\eta_t, \mathbf{e}_{g,t}) \sim \mathcal{N}(0, \mathbf{I})$. We set $(\mathbf{u}_t, \mathbf{e}_{g,t}, \eta_t)$ be serially independent, but still allow the same cross-sectional dependence among \mathbf{u}_t . This allows us to focus on the effect of over-estimating factors. The r -dimensional \mathbf{f}_t are independent standard normal. We set the true $\beta = 1$, $\theta = \mathbf{v} = (1, -1.5, 0.5, 0, \dots, 0)$ and $T = N = 200$.

We employ the diversified factor estimator described in Section 3.2 with various working number of factors R , and compare with the benchmark “double-selection” method of Belloni, Chernozhukov, and Hansen (2014). In particular, we consider two settings:

- (i) $r = 0$: there are no factors so \mathbf{x}_t itself is weakly dependent.
- (ii) $r = 2$: there are two factors driving \mathbf{x}_t . Set $\alpha = 1$ so both factors are strong.

We calculate the standardized estimates: $N := \hat{\sigma}_{\eta, \mathbf{g}}^{-1} \hat{\sigma}_{\mathbf{g}}^2 \sqrt{T}(\hat{\beta} - \beta)$, where the standard error is the estimated feasible one. Our theory shows that the sampling distribution of z should be approximately standard normal.

Figures 2 and 3 plot the histograms of the standardized estimates over 200 replications, superimposed with the standard normal density. The histogram is scaled to be a density function. We present the results when the initial transformation are used as weights for the diversified factors. The results from characteristics and Hadamard weights are very similar. When $r = 0$, while it is expected that the double selection performs very well, as is shown in Figure 3, using $R \geq 1$ factors also produces z -statistics whose distribution is also close to the standard normality. This shows that the factor-augmented method is robust to the absence of factor structures. On the other hand, when $r = 2$, the factor-augmented method continues to perform well. In contrast, the double selection is severely biased, and the distribution of its z -statistic is far off from the standard normality.

The first three panels employ the diversified factor estimator with R number of working factors. The last panel uses the double

selection, which directly selects among \mathbf{x}_t . The weights used are the initial transformations ($t = 0$) so that the i th row of \mathbf{W} is $(x_{i,0}, x_{i,0}^2, \dots, x_{i,0}^R)$ for $i \leq N$.

5.4. Testing the Specification of Empirical Factors

In the last simulation study, we study the size and power of the test statistic for $H_0 : \mathbf{g}_t = \theta \mathbf{f}_t$ for some $r \times r$ invertible matrix θ . Here, \mathbf{g}_t is a vector of known “empirical factors” that applied researchers propose to approximate the true factors. We generate

$$\mathbf{g}_t = \theta \mathbf{f}_t + \gamma \mathbf{h}_t, \quad t \leq T,$$

where θ is an r -dimensional identity matrix, and $(\mathbf{f}_t, \mathbf{h}_t) \sim \mathcal{N}(0, \mathbf{I})$. Here γ governs the strength of the alternatives. We assume that \mathbf{u}_t is serially independent normal generated from $\mathcal{N}(0, \Sigma_N)$, with Σ_N as in (21), pertaining the same cross-sectional dependence. We set $R = r = 2$ and fix $N = 200$. In each of the simulations, we calculate the test statistic as defined in Section 3.4, and set the significance level to 0.05. We use the SCAD-thresholding to estimate Σ_u for both $\widehat{\text{MEAN}}$ and $\widehat{\sigma}$.

Table 2 presents the rejection probability over 1000 replications, with $\gamma = 0$ representing the size of the test. Above all, the results look satisfactory with controlled size and reasonable powers, while weights using initial transformations have some size distortions.

6. Conclusion

We propose simple estimators of the latent factors using cross-sectional projections of the panel data, by weighted averages. These weights are chosen to diversify away the idiosyncratic components, resulting in “diversified factors.” Because the projections are conducted cross-sectionally, they are robust to serial conditions, easy to analyze due to data-independent weights, and work even for finite length of time series. We formally prove

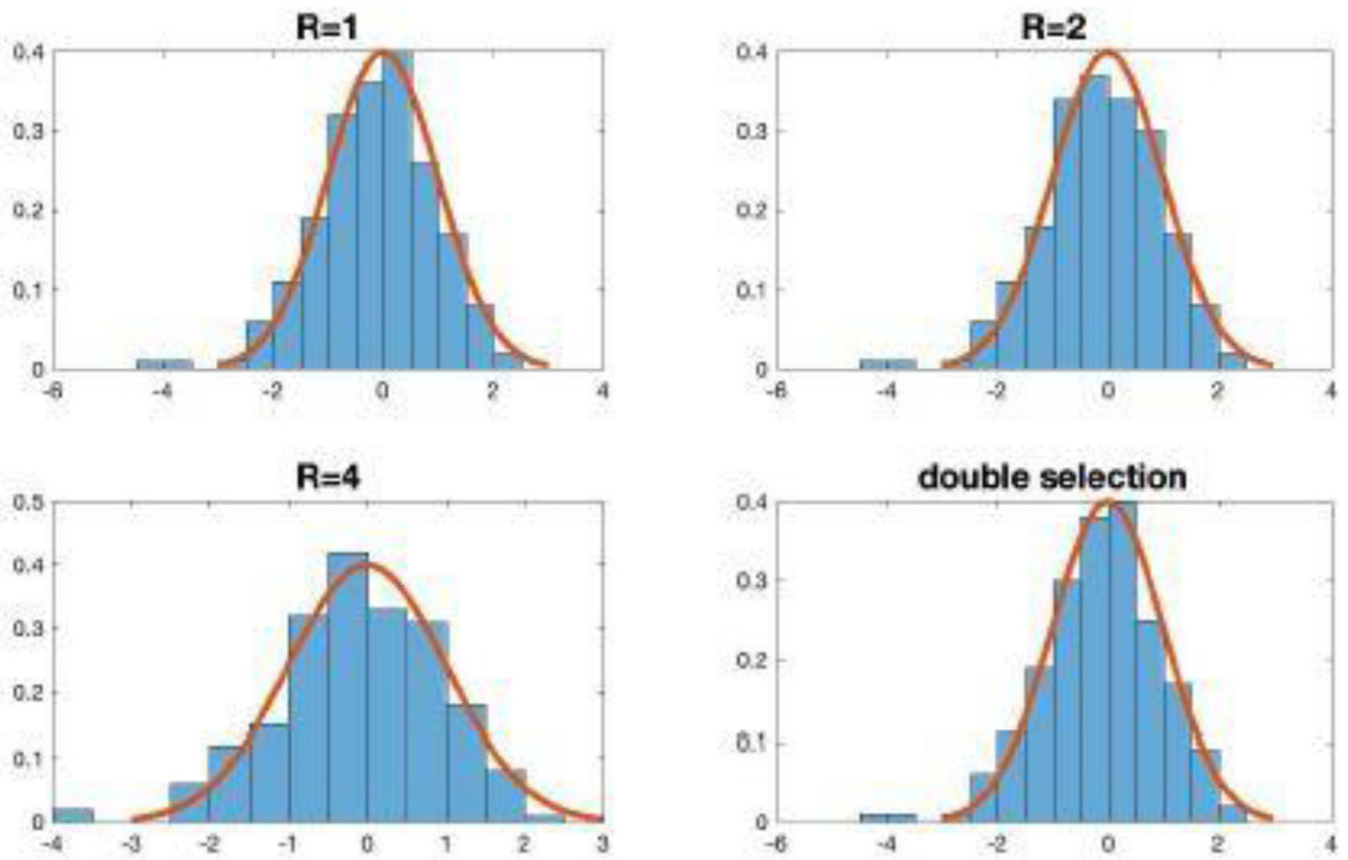
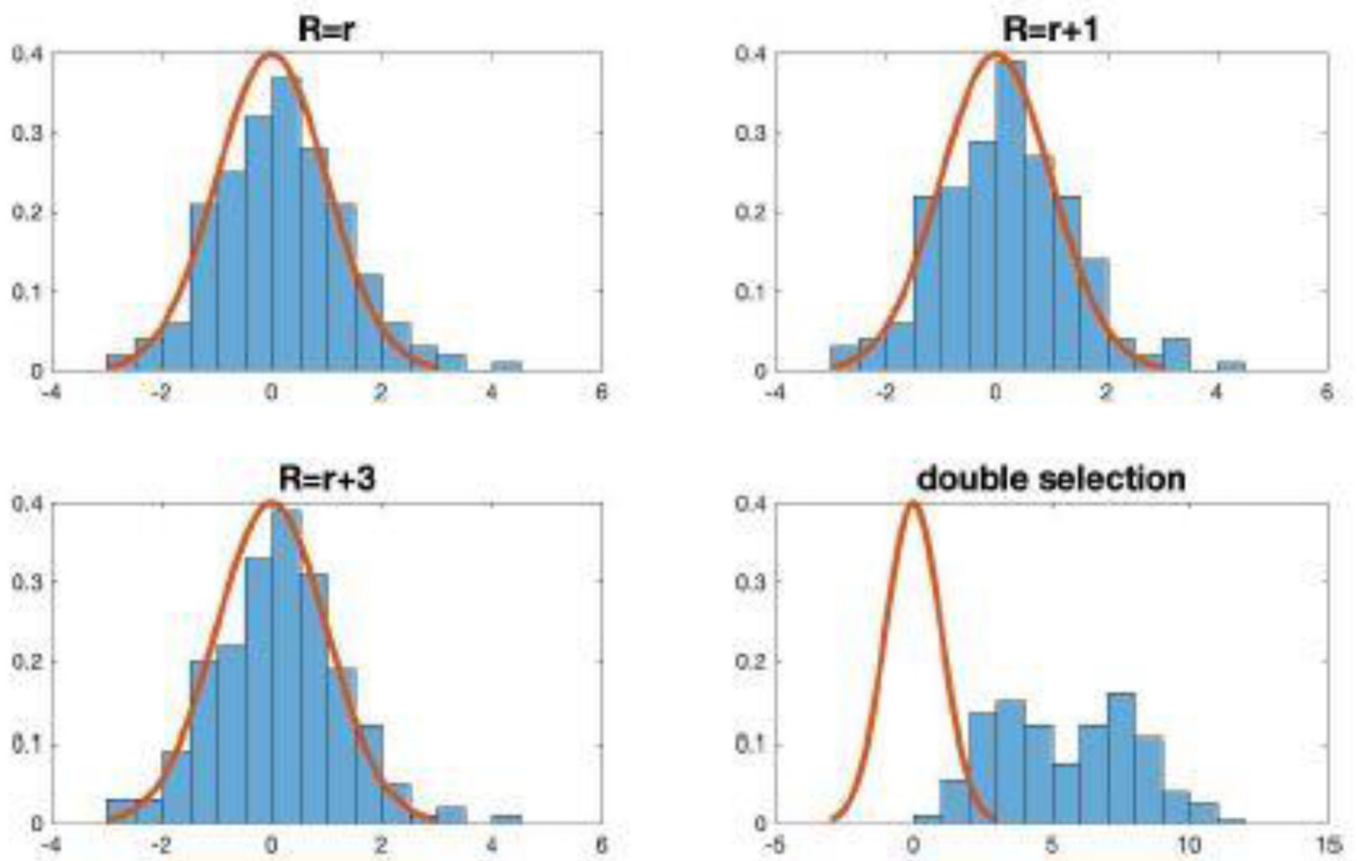
Figure 2. True $r = 0$.Figure 3. True $r = 2$.

Table 2. Probability of rejection at level 0.05.

γ	T	Characteristic weights	Hadamard weights	Initial transformation
0	100	0.054	0.046	0.065
	200	0.052	0.047	0.074
0.2	100	1.000	0.998	1.000
	200	0.975	1.000	1.000

NOTE: γ represents the strength of alternatives.

that this procedure is robust to over-estimating the number of factors, and illustrate it in several applications. We also recommend several choices for the diversified weights.

Supplementary Materials

The supplementary material contains all the technical proofs.

Funding

Jianqing Fan's research is supported by NSF grants DMS-1662139 and DMS-1712591.

References

- Ahn, S., and Horenstein, A. (2013), "Eigenvalue Ratio Test for the Number of Factors," *Econometrica*, 81, 1203–1227. [910,912,913]
- Antoniadis, A., and Fan, J. (2001), "Regularized Wavelet Approximations," *Journal of the American Statistical Association*, 96, 939–967. [916]
- Bai, J. (2003), "Inferential Theory for Factor Models of Large Dimensions," *Econometrica*, 71, 135–171. [909]
- Bai, J., and Li, K. (2012), "Statistical Analysis of Factor Models of High Dimension," *The Annals of Statistics*, 40, 436–465. [916]
- Bai, J., and Liao, Y. (2017), "Inferences in Panel Data With Interactive Effects Using Large Covariance Matrices," *Journal of Econometrics*, 200, 59–78. [917]
- Bai, J., and Ng, S. (2002), "Determining the Number of Factors in Approximate Factor Models," *Econometrica*, 70, 191–221. [910]
- (2006), "Confidence Intervals for Diffusion Index Forecasts and Inference for Factor-Augmented Regressions," *Econometrica*, 74, 1133–1150. [913]
- Barigozzi, M., and Cho, H. (2018), "Consistent Estimation of High-Dimensional Factor Models When the Factor Number Is Over-Estimated," arXiv no. 1811.00306. [910,918,919]
- Belloni, A., Chernozhukov, V., and Hansen, C. (2014), "Inference on Treatment Effects After Selection Among High-Dimensional Controls," *The Review of Economic Studies*, 81, 608–650. [914,915,916,922]
- Bickel, P., and Levina, E. (2008), "Covariance Regularization by Thresholding," *The Annals of Statistics*, 36, 2577–2604. [916]
- Bickel, P., Ritov, Y., and Tsybakov, A. (2009), "Simultaneous Analysis of Lasso and Dantzig Selector," *The Annals of Statistics*, 37, 1705–1732. [915]
- Carhart, M. M. (1997), "On Persistence in Mutual Fund Performance," *Journal of Finance*, 52, 57–82. [916]
- Chudik, A., Pesaran, M. H., and Tosetti, E. (2011), "Weak and Strong Cross-Section Dependence and Estimation of Large Panels," *The Econometrics Journal*, 14, C45–C90. [910,912]
- Connor, G., and Korajczyk, R. A. (1986), "Performance Measurement With the Arbitrage Pricing Theory: A New Framework for Analysis," *Journal of Financial Economics*, 15, 373–394. [909]
- Connor, G., Matthias, H., and Linton, O. (2012), "Efficient Semiparametric Estimation of the Fama–French Model and Extensions," *Econometrica*, 80, 713–754. [910,918]
- Doz, C., Giannone, D., and Reichlin, L. (2011), "A Two-Step Estimator for Large Approximate Dynamic Factor Models Based on Kalman Filtering," *Journal of Econometrics*, 164, 188–205. [920,922]
- Fama, E. F., and French, K. R. (1992), "The Cross-Section of Expected Stock Returns," *Journal of Finance*, 47, 427–465. [916]
- Fan, J., Ke, Y., Sun, Q., and Zhou, W.-X. (2019), "Farmtest: Factor-Adjusted Robust Multiple Testing With Approximate False Discovery Control," *Journal of the American Statistical Association*, 114, 1880–1893. [918]
- Fan, J., Ke, Y., and Wang, K. (2020), "Factor-Adjusted Regularized Model Selection," *Journal of Econometrics*, 216, 71–85. [914,915]
- Fan, J., Liao, Y., and Mincheva, M. (2013), "Large Covariance Estimation by Thresholding Principal Orthogonal Complements" (with discussion), *Journal of the Royal Statistical Society, Series B*, 75, 603–680. [916,920]
- Fan, J., Liao, Y., and Wang, W. (2016), "Projected Principal Component Analysis in Factor Models," *The Annals of Statistics*, 44, 219–254. [910,918,919]
- Fan, J., Liao, Y., and Yao, J. (2015), "Power Enhancement in High Dimensional Cross-Sectional Tests," *Econometrica*, 83, 1497–1541. [916,917]
- Forni, M., Hallin, M., Lippi, M., and Reichlin, L. (2005), "The Generalized Dynamic Factor Model: One-Sided Estimation and Forecasting," *Journal of the American Statistical Association*, 100, 830–840. [920,922]
- Gagliardini, P., Ossola, E., and Scaillet, O. (2016), "Time-Varying Risk Premium in Large Cross-Sectional Equity Data Sets," *Econometrica*, 84, 985–1046. [918]
- Goncalves, S., and Perron, B. (2020), "Bootstrapping Factor Models With Cross Sectional Dependence," *Journal of Econometrics*, 218, 476–495. [916]
- Hallin, M., and Liška, R. (2007), "Determining the Number of Factors in the General Dynamic Factor Model," *Journal of the American Statistical Association*, 102, 603–617. [910,920,922]
- Hansen, C., and Liao, Y. (2018), "The Factor-Lasso and k -Step Bootstrap Approach for Inference in High-Dimensional Economic Applications," *Econometric Theory*, 35, 465–509. [914]
- Johnstone, I. M., and Lu, A. Y. (2009), "On Consistency and Sparsity for Principal Components Analysis in High Dimensions," *Journal of the American Statistical Association*, 104, 682–693. [909]
- Juodis, A., and Sarafidis, V. (2020), "A Linear Estimator for Factor-Augmented Fixed-t Panels With Endogenous Regressors," Tech. Rep., Monash University, Department of Econometrics and Business Statistics. [919]
- Karabiyik, H., Reese, S., and Westerlund, J. (2017), "On the Role of the Rank Condition in CCE Estimation of Factor-Augmented Panel Regressions," *Journal of Econometrics*, 197, 60–64. [912]
- Karabiyik, H., Urbain, J.-P., and Westerlund, J. (2019), "CCE Estimation of Factor-Augmented Regression Models With More Factors Than Observables," *Journal of Applied Econometrics*, 34, 268–284. [910,912]
- Kim, S., Korajczyk, R. A., and Neuhierl, A. (2018), "Arbitrage Portfolios in Large Panels," available at SSRN. [919]
- Li, G., Yang, D., Nobel, A. B., and Shen, H. (2016), "Supervised Singular Value Decomposition and Its Asymptotic Properties," *Journal of Multivariate Analysis*, 146, 7–17. [918]
- Li, H., Li, Q., and Shi, Y. (2017), "Determining the Number of Factors When the Number of Factors Can Increase With Sample Size," *Journal of Econometrics*, 197, 76–86. [910]
- Ludvigson, S., and Ng, S. (2007), "The Empirical Risk–Return Relation: A Factor Analysis Approach," *Journal of Financial Economics*, 83, 171–222. [913]
- Moon, R., and Weidner, M. (2015), "Linear Regression for Panel With Unknown Number of Factors as Interactive Fixed Effects," *Econometrica*, 83, 1543–1579. [910]
- Park, B. U., Mammen, E., Härdle, W., and Borak, S. (2009), "Time Series Modelling With Semiparametric Factor Dynamics," *Journal of the American Statistical Association*, 104, 284–298. [910]
- Pesaran, H. (2006), "Estimation and Inference in Large Heterogeneous Panels With a Multifactor Error Structure," *Econometrica*, 74, 967–1012. [910,912]
- Robinson, P. M. (1988), "Root- n -Consistent Semiparametric Regression," *Econometrica*, 56, 931–954. [914]
- Stock, J., and Watson, M. (2002), "Forecasting Using Principal Components From a Large Number of Predictors," *Journal of the American Statistical Association*, 97, 1167–1179. [909,913]
- Wang, J., Zhao, Q., Hastie, T., and Owen, A. B. (2017), "Confounder Adjustment in Multiple Hypothesis Testing," *The Annals of Statistics*, 45, 1863. [918]
- Westerlund, J., and Urbain, J.-P. (2015), "Cross-Sectional Averages Versus Principal Components," *Journal of Econometrics*, 185, 372–377. [910]