

Fan, Jianqing; Masini, Ricardo; Medeiros, Marcelo C.

Working Paper

Bridging factor and sparse models

Texto para discussão, No. 681

Provided in Cooperation with:

Departamento de Economia, Pontifícia Universidade Católica do Rio de Janeiro

Suggested Citation: Fan, Jianqing; Masini, Ricardo; Medeiros, Marcelo C. (2021) : Bridging factor and sparse models, Texto para discussão, No. 681, Pontifícia Universidade Católica do Rio de Janeiro (PUC-Rio), Departamento de Economia, Rio de Janeiro

This Version is available at:

<https://hdl.handle.net/10419/249729>

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.

TEXTO PARA DISCUSSÃO

No. 681

Bridging factor and sparse models

Jianqing Fan
Ricardo P. Masini
Marcelo C. Medeiros



Bridging factor and sparse models

Jianqing Fan

Department of Operations Research and Financial Engineering
Princeton University
E-mail: jqfan@princeton.edu

Ricardo Masini

Center for Statistics and Machine Learning, Princeton University
E-mail: rmasini@princeton.edu

Marcelo C. Medeiros

Department of Economics, Pontifical Catholic University of Rio de Janeiro
E-mail: mcm@econ.puc-rio.br

March 5, 2021

Abstract

Factor and sparse models are two widely used methods to impose a low-dimensional structure in high dimension. They are seemingly mutually exclusive. In this paper, we propose a simple lifting method that combines the merits of these two models in a supervised learning methodology that allows to efficiently explore all the information in high-dimensional datasets. The method is based on a flexible model for panel data, called factor-augmented regression model with both observable, latent common factors, as well as idiosyncratic components as high-dimensional covariate variables. This model not only includes both factor regression and sparse regression as specific models but also significantly weakens the cross-sectional dependence and hence facilitates model selection and interpretability. The methodology consists of three steps. At each step, remaining cross-section dependence can be inferred by a novel test for covariance structure in high-dimensions. We developed asymptotic theory for the factor-augmented sparse regression model and demonstrated the validity of the multiplier bootstrap for testing high-dimensional covariance structure. This is further extended to testing high-dimensional partial covariance structures. The theory and methods are further supported by an extensive simulation study and applications to the construction of a partial covariance network of the financial returns and a prediction exercise for a large panel of macroeconomic time series from FRED-MD database.

JEL Codes: C22, C23, C32, C33.

Keywords: Factor models, sparse regression, high-dimensional, supervised learning, hypothesis testing, covariance structure, partial covariance structure.

Acknowledgments: Medeiros gratefully acknowledges the partial financial support from CNPq and CAPES. We are very grateful to Alexander Giessing, Bruno Ferman, Caio Almeida, Claudio Flores, Conrado Garcia, Gilberto Boareto, Gustavo Bulhões, Henrique Pires, Marcelo Fernandes, Marcelo J. Moreira, Nathalie Gimenes, Rodrigo Targino, and Yuri Saporito for helpful discussions and comments. We also thank the participants during the Wilks Seminar at Princeton University for valuable comments. Finally, we are deeply grateful to Eduardo F. Mendes and Michael Wolf for the careful reading of the paper and the many insightful discussions which led to a much improved version of this manuscript.

1 Introduction

With the emergence of new and large datasets, the correct characterization of the dependence among variables is of substantial importance. Usually, to achieve this goal, the literature has followed two seemingly orthogonal tracks. On the one hand, factor models have become an essential tool to summarize information in large datasets under the assumption that the remaining dependence structure is negligible. For instance, panel factor models are applied now to a wide variety of important applications, ranging from forecasting (macroeconomic) variables and asset pricing models to causal inference in applied microeconomics and network analysis. On the other hand, there have been major advances on parameter estimation in ultra high-dimensions under the assumption of sparsity or weak-sparsity. That is, a variable depends only on a (very) small subset of the other variables.

In this paper, we take an alternative route and combine the best of the two worlds described above in order to better characterize the dependence structure of high-dimensional data. More specifically, we consider that the covariance structure of a large set of variables, organized in a panel data format, is characterized as a combination of a factor structure, where factors can be either observed, unobserved, or both, and a weakly-sparse idiosyncratic component. This formulation is general enough in order to accommodate a very large number of data generating processes of interest in economics, finance, and related areas. The proposed methodology has two ingredients: a three-step estimation procedure and a new test for structure in high dimensional (partial) covariance matrices. The steps of the estimation procedure are as follows. In the first one, we take the original data and remove the effects of any observed factors. These factors can be deterministic terms such as seasonal dummies and/or trends or any other observed covariates. The first step can be parametric or nonparametric, low or high dimensional. A latent factor model is then estimated using the residuals from the first stage. Finally, in a final step we model the dependence among idiosyncratic terms as a weakly sparse regression estimated by the Least Absolute Shrinkage and Selection Operator (LASSO). At each step, the null-hypothesis of no remaining cross-section dependence can be tested by the proposed test for the (partial) covariance structure in high-dimensions.

1.1 Motivation

Let $\mathbf{Y}_t := (Y_{1t}, \dots, Y_{nt})'$ be a random vector generated by a factor model as $Y_{it} = \boldsymbol{\lambda}'_i \mathbf{F}_t + U_{it}$, for $i = 1, \dots, n$, $t = 1, \dots, T$, where $\boldsymbol{\Sigma} := \mathbb{E}(\mathbf{U}_t \mathbf{U}'_t)$, with $\mathbf{U}_t := (U_{1t}, \dots, U_{nt})'$, is not necessarily diagonal. Fix one component of interest $i \in \{1, \dots, n\}$, which serve as a response variable. Consider

the following prediction models:

$$\mathcal{M}_1 : \mathbb{E}(Y_{it}|\mathbf{Y}_{-it}), \quad \mathcal{M}_2 : \mathbb{E}(Y_{it}|\mathbf{F}_t), \quad \text{and} \quad \mathcal{M}_3 : \mathbb{E}(Y_{it}|\mathbf{F}_t, \mathbf{U}_{-it}), \quad (1.1)$$

where \mathbf{Y}_{-it} and \mathbf{U}_{-it} are, respectively, vectors with the elements of \mathbf{Y}_t and \mathbf{U}_t excluding the i -th entry. Note that model \mathcal{M}_3 is indeed the factor augmented regression model since it is the same as $\mathbb{E}(Y_{it}|\mathbf{F}_t, \mathbf{Y}_{-it})$. As the paper will mainly focus on linear regressions, we will refer more specifically $\widetilde{\mathcal{M}}_3$ below as the factor-augmented regression model.

Suppose that we observe both \mathbf{F}_t and \mathbf{U}_{-it} . Which one of three models above is best in terms of mean square error (MSE) for prediction? Comparison between \mathcal{M}_1 and \mathcal{M}_2 is not clear since it depends, among others, on the magnitude of Σ relative to $\Lambda'\Lambda$, where $\Lambda := (\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_n)'$. However, since the σ -algebras generated by \mathbf{Y}_{-it} and \mathbf{F}_t are both included in the σ -algebra generated by $(\mathbf{F}_t, \mathbf{U}_{-it})$, it is not surprising that $\text{MSE}(\mathcal{M}_3) \leq \min[\text{MSE}(\mathcal{M}_1), \text{MSE}(\mathcal{M}_2)]$. The same will hold true if we replace the models in (1.1) by their best linear projections, which we denote by $\widetilde{\mathcal{M}}_j$ for $j \in \{1, 2, 3\}$, since the linear space $\widetilde{\mathcal{M}}_3$ is the largest. In this case, we can explicitly write the “gains” of $\widetilde{\mathcal{M}}_3$ when compared to $\widetilde{\mathcal{M}}_1$ and $\widetilde{\mathcal{M}}_2$:

$$\begin{aligned} \text{MSE}(\widetilde{\mathcal{M}}_3) - \text{MSE}(\widetilde{\mathcal{M}}_1) &= -\boldsymbol{\theta}'_i \boldsymbol{\Sigma}_{-i,-i} \boldsymbol{\theta}_i \\ \text{MSE}(\widetilde{\mathcal{M}}_3) - \text{MSE}(\widetilde{\mathcal{M}}_2) &= -\boldsymbol{\Delta}_{1i} \boldsymbol{\Delta}'_{1i} - \boldsymbol{\Delta}'_{2i} \boldsymbol{\Sigma}_{-i,-i} \boldsymbol{\Delta}_{2i}, \end{aligned}$$

where $\boldsymbol{\theta}_i$ and $\boldsymbol{\beta}_i$ are the coefficients of the projection of U_{it} onto \mathbf{U}_{-it} and the coefficients of the projection of X_{it} onto \mathbf{X}_{-it} , respectively; $\boldsymbol{\Sigma}_{-i,-i}$ is Σ excluding the i -th row and column; $\boldsymbol{\Delta}_{1i} := \boldsymbol{\Lambda}_i - \boldsymbol{\beta}'_i \boldsymbol{\Lambda}_{-i}$ and $\boldsymbol{\Delta}_{2i} := \boldsymbol{\beta}_i - \boldsymbol{\theta}_i$. From the previous expressions, it becomes evident that both $\widetilde{\mathcal{M}}_1$ and $\widetilde{\mathcal{M}}_2$ are restrictions on $\widetilde{\mathcal{M}}_3$. Broadly speaking, whenever one does not expect to have an *exact* factor model, there are potential gains of taking into account the contribution of the idiosyncratic components \mathbf{U}_{-it} . Therefore, we use $\widetilde{\mathcal{M}}_3$ as the base model for the estimation methodology described in Section 2.2.

1.2 Main Contributions and Comparison with the Literature

The contributions of this paper are multi-fold. First, our methodology bridges the gap between two apparently competing methods for high-dimensional modeling; see, for example the discussion in Giannone et al. (2018) and Fan et al. (2020). This yields a vast number of potential applications and

spin-offs. For instance, in Fan et al. (2020), we apply the methods developed in here to evaluate the effects of interventions and we contribute to the literature on synthetic controls and related methods by combining the approaches of Gobillon and Magnac (2016) and Carvalho et al. (2018). Therefore, in our setup both a common factor structure and weak sparsity can coexist.¹

Second, our results can also serve as a diagnostic and misspecification tool. For panel data models with interactive fixed effects as in Moon and Weidner (2015) and Bai and Liao (2017), our test can be directly applied to uncover the dependence structure among cross-sectional units before and after accounting for common factor components. If the factor structure is informative enough, we expect the idiosyncratic covariance matrix to be almost sparse. If this is not the case, we may have possibly underestimated the number of factors. One popular application is in asset pricing as discussed in Gagliardini et al. (2019) and in the empirical section of this paper. There are a huge number of proposed factors as described in Feng et al. (2020), Giglio and Xiu (2020), and Gu et al. (2020). We can apply our methodology not only to test for omitted factors, but, as well, to estimate network connections among firms as in Diebold and Yilmaz (2014) and Brownlees et al. (2020). Finally, as a diagnostic tool, our paper tackle the same problem as Gagliardini et al. (2019). However, we take an alternative solution strategy which relies on a much different set of hypothesis; see also Gagliardini et al. (2020).

Third, the methodology proposed here contributes to the forecasting literature. For instance, in the second application considered in this paper, we build forecasting models for a large cross-section of macroeconomic variables. We call this method the **FarmPredict**. We show that the combination of factors and a sparse regression strongly outperforms the traditional principal component regression as in Stock and Watson (2002a,b). Therefore, **FarmPredict** can be an additional contribution to the forecasting and machine learning toolkit. The method can be easily extended to a multivariate setting combining factor-augmented vector autoregressions (FAVAR) as in Bernanke et al. (2005) with sparse vector models as in Kock and Callot (2015) and Masini et al. (2019).

Fourth, we show consistency of factor estimation based on the residuals of a first-step regression. Our results hold for both parametric (linear or nonlinear) and nonparametric first stage. A high-dimensional first stage is also allowed. Note that, current results in the literature consider that factors are estimated based on observed data and our derivations favor a much more flexible and general setup (Bai and Ng, 2002, 2003, 2006). More specifically, our methodology favors settings

¹Sparsity and factor models can also coexist in the framework of sparse principal components; see Fan et al. (2020).

where there are both observed and latent factors, as well as trend-stationary data. In the later, the trend can be first removed by (nonparametric) first-stage regression.

Fifth, we also contribute to the LASSO literature. LASSO can not be model selection consistent for highly correlated variables. Through the decomposition of covariates into factors and idiosyncratic components, we decorrelate the variables and make the model selection condition much easier to hold; see, for example, (Fan et al., 2020). We show consistency of the estimates based on residuals of the previous steps. Our results are derived under restrictions on the population covariance matrix of the data and not on the estimated one, as it is usual in many papers. See, for example, van de Geer and Bühlmann (2009). Furthermore, we derive our results under much mild conditions that the ones considered in (Fan et al., 2020).

Finally, we extend the results in Chernozhukov et al. (2013, 2018) to strong-mixing data in order to construct hypothesis tests for covariance and partial covariance structure in high dimensions.² This step is necessary for econometrics and financial applications. As side results, in order to develop the test we first show consistency of kernel-based estimation of a high-dimensional long-run covariance matrix of dependent process. This is a new result with important consequences for the theory of high-dimensional regression with dependent errors. We also establish consistency of a new estimator of the partial covariance matrix in high-dimensions and strong-mixing data. Our proposed tests can be used to infer, for instance, if the (partial) covariance matrix of a high-dimensional random vector is diagonal or block-diagonal. More generally, we can test any pre-defined structure. Furthermore, we show that the test remains valid when we use the residuals from a previous step estimation to compute the covariance matrix. This result allows us to to apply the test to the three-stage estimation procedure proposed in this paper. Although our results are derived under the assumption that the number of factors is known, simulation results presented in the paper provides evidence that the test have good finite-sample properties even when the number of factors is determined by data-driven methods commonly found in the literature. Over the past years, a vast number of papers proposed different methods to test for covariance structure in high dimensions. See, for example, Ledoit and Wolf (2002), Chen et al. (2010), Onatski et al. (2013), Cai and Ma (2013), Li and Qin (2014), Zheng et al. (2019), Cai et al. (2016), Zheng et al. (2019), and Guo and Tang (2020), among many others.³ To the best of our knowledge, we complement all the previous papers by simultaneously considering high-dimensions, strong-mixing data with mild distributional

²Recently, Giessing and Fan (2020) also extended the results in Chernozhukov et al. (2013). However, their setup is very different from ours and the authors only consider the case of independent and identically distributed data.

³For a nice recent review, see Cai (2017).

assumptions, and pre-estimation when constructing tests for both covariance and partial covariance structure.

1.3 Organization of the Paper

In addition to this Introduction, the paper is organized as follows. We present the model setup and assumptions in Section 2. The theoretical results are presented in Section 3 with practical guides given in Section 4. We depict the results of a simulation experiment in Section 5 and discuss the empirical application in Section 6. Section 7 concludes. All proofs are deferred to the Appendix.

1.4 Notation

All random variables (real-valued scalars, vectors and matrices) are defined in a common probability space $(\Omega, \mathcal{F}, \mathbb{P})$. We denote random variables by an upper case letter, X for instance, and its realization by a lower case letter, $X = x$. The expected value operator is with respect to the \mathbb{P} law such that $\mathbb{E}X := \int_{\Omega} X(\omega) d\mathbb{P}(\omega)$. Matrices and vectors are written in bold letters \mathbf{X} . Except for the number of factors, r , and number of covariates, k , defined below, all other dimensions are allowed to depend on the sample size (T). However, we omit this dependency throughout the paper to avoid clustering the notation prematurely.

We use $\|\cdot\|_p$ to denote the ℓ^p norm for $p \in [1, \infty]$, such that for a d -dimensional (possibly random) vector $\mathbf{X} = (X_1, \dots, X_d)'$, we have $\|\mathbf{X}\|_p := (\sum_{i=1}^d |X_i|^p)^{1/p}$ for $p \in [1, \infty)$ and $\|\mathbf{X}\|_{\infty} := \sup_{i \leq d} |X_i|$. If \mathbf{X} is a $(m \times n)$ possibly random matrix then $\|\mathbf{X}\|_p$ denotes the matrix ℓ^p -induced norm and $\|\mathbf{X}\|_{\max}$ denotes the maximum entry in absolute terms of the matrix \mathbf{X} . Note that whenever \mathbf{X} is random, then $\|\mathbf{X}\|_p$ for $p \in [1, \infty]$ and $\|\mathbf{X}\|_{\max}$ are random variables. We also reserve the symbol $\|\cdot\|$ without subscript for the Euclidean norm $\|\cdot\| := \|\cdot\|_2$ for both vectors and matrices.

For any convex function $\psi : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ such that $\psi(0) = 0$ and $\psi(x) \rightarrow \infty$ as $x \rightarrow \infty$ and (real-valued) random variable X , we denote its Orlicz norm by $\|X\|_{\psi}$, which is defined by $\|X\|_{\psi} := \inf \left\{ C > 0 : \mathbb{E} \left[\psi \left(\frac{|X|}{C} \right) \right] \leq 1 \right\}$. Since we are only concerned with polynomial and exponential tails we restrict ourselves to Orlicz norm induces by the class of function defined by (3.3). Evidently, as opposed to $\|X\|_p$, $\|X\|_{\psi_p}$ is always a non-negative non-random scalar. We do not abide to any convention to apply Orlicz norm to vector or matrices to avoid confusion.

For any vector \mathbf{X} , $\text{diag}(\mathbf{X})$ denote the diagonal matrix whose diagonal is the elements of \mathbf{X} . $\mathbb{1}(A)$ is an indicator function on the event A , i.e, $\mathbb{1}(A) = 1$ if A is true and 0 otherwise. We adopt

the Landau big/small O , o notation and the “in probability” O_P and o_p analogues. We say that x is of the same order of y , $x \asymp y$, if both $x = O(y)$ and $y = O(x)$. We write $X \asymp_P Y$ if both $X = O_P(Y)$ and $Y = O_P(X)$. Unless stated otherwise, the asymptotics are taken as $T \rightarrow \infty$, where T is the time-series dimension, and the $o(1)$ and $o_P(1)$ are with respect to the limit as $T \rightarrow \infty$. We denote convergence in probability and in distribution by “ \xrightarrow{P} ” and “ \Rightarrow ”, respectively.

2 Setup and Method

2.1 Data Generating Process

We apply the test for three-stage estimation procedure for a very general panel data model, which is rich enough in order to nest several important cases in economics, finance and related areas. More specifically, we define the following the Data Generating Process (DGP).

Assumption 1 (DGP). *The process $\{Y_{it} : 1 \leq i \leq n, t \geq 1\}$ is generated by*

$$Y_{it} = \gamma_i' \mathbf{X}_{it} + \underbrace{\boldsymbol{\lambda}_i' \mathbf{F}_t}_{=: R_{it}} + U_{it} \quad (2.1)$$

where \mathbf{X}_{it} is a k -dimensional observable (random) vector which may also include a constant term, \mathbf{F}_t is a r -dimensional vector of common latent factors, and U_{it} is a zero mean idiosyncratic shock.⁴ The unknown parameters are $\gamma_i \in \mathbb{R}^k$, the factor loadings $\boldsymbol{\lambda}_i$, and the covariance matrix of the idiosyncratic shocks. Finally, we assume that \mathbf{X}_{it} , \mathbf{F}_t and U_{it} are mutually uncorrelated.

Remark 1. *In Assumption 1 we consider that k , the dimension of \mathbf{X}_{it} is finite and fixed. Furthermore, the relation between Y_{it} and \mathbf{X}_{it} is linear. This is for the sake of exposition. However, the theoretical results in this paper are written in terms of the consistency rate of the first-step estimation. Therefore, the DGP can be made much more general by just changing the rates.*

Example 1 (Asset Pricing Models). *Suppose Y_{it} is the return of an asset i at time t and let $\mathbf{X}_{it} := \mathbf{X}_t$ be a set of k observable risk factors, such as the market returns and or Fama-French factors as in, for example, Fama and French (1993,2015). \mathbf{F}_t can be a set of additional, non observable, risk factors. Several asset pricing models, such as the Capital Asset Pricing Model (CAPM) of the Arbitrage Pricing Theory (APT) model, are nested into this general framework.*

⁴For simplicity, we assume that all the units i have the same number of covariates (k). The framework can certainly accommodate situations where k_i is a function of i .

Example 2 (Networks). *Model (2.1) also complements the network specifications discussed in Barigozzi and Hallin (2016,2017b) and Barigozzi and Brownlees (2019). Furthermore, the test proposed here can be used to detect networks links as in Diebold and Yilmaz (2014) and Brownlees et al. (2020). For example, Y_{it} can be the (realized) volatility of financial assets and $\mathbf{X}_{it} := \mathbf{X}_t$ can be volatility factors as in Brito et al. (2018) and Andreou and Ghysels (2021).*

Example 3 (Panel Data Models with Iterative Fixed-Effects). *Model (2.1) is the panel model with iterative fixed-effects considered in Gobillon and Magnac (2016), where the authors propose an alternative to the Synthetic Control method of Abadie and Gardeazabal (2003) and Abadie et al. (2010) to evaluate the effects of regional policies. Model (2.1) is also in the heart of the FarmTreat method of Fan et al. (2020) and the model discussed in Moon and Weidner (2015).*

Example 4 (FAVAR). *In the case where the index i represents a different dependent (endogenous) variable and U_{it} is a dependent process, model (2.1) turns out to be equivalent to the Factor Augmented Vector Autoregressive (FAVAR) model of Bernanke et al. (2005). In this case, \mathbf{X}_{it} may also include lagged dependent variables.*

2.2 Three-Stage Method

The method proposed here for estimation, inference and prediction consists of three stages where at the end of each stage, the covariance structure of the residuals is tested.

1. For each $i \in \{1, \dots, n\}$ run the regression:

$$Y_{it} = \boldsymbol{\gamma}'_i \mathbf{X}_{it} + R_{it}, \quad t \in \{1, \dots, T\},$$

and compute $\hat{R}_{it} := Y_{it} - \hat{\boldsymbol{\gamma}}'_i \mathbf{X}_{it}$. The first stage may consist of a regression on a constant, a deterministic time trend and seasonal dummies, for instance, or, as in Example 1, a regression on observed factors. After removing the contribution from the observables, we can use the test for the null hypothesis of no remaining (partial) covariance structure to check if the (partial) covariance of R_{it} is dense or sparse. If it is dense we move to Step 2. Otherwise, we jump directly to Step 3. This first parametric, low dimensional step can be replaced by a nonlinear/nonparametric regression or by a high-dimensional model, when, for example, the number of observed factors is large. This will be discussed more in the subsequent sections.

2. Write $\mathbf{R}_t := (R_{1t}, \dots, R_{nt})'$ and $\mathbf{R}_t = \mathbf{\Lambda}\mathbf{F}_t + \mathbf{U}_t$. The second step consists of estimating $\mathbf{\Lambda}$ and \mathbf{F}_t for $t = 1, \dots, T$ using $\hat{\mathbf{R}}_t$ through principal component analysis (PCA) and compute

$$\hat{\mathbf{U}}_t = \hat{\mathbf{R}}_t - \hat{\mathbf{\Lambda}}\hat{\mathbf{F}}_t.$$

After estimating the factors and loadings, we apply our testing procedure to test for remaining covariance structure in \mathbf{U}_t . The second-step estimation can be carried out also by dynamic factor models as in Barigozzi and Hallin (2016,2017,2020) or Barigozzi et al. (2020).

3. Now, define $\hat{\mathbf{U}}_{-it} := (\hat{U}_{1t}, \dots, \hat{U}_{i-1,t}, \hat{U}_{i+1,t} \dots \hat{U}_{nt})'$. The third estimation step consists of a sparse regression to estimate the following model for each $i \in \{1, \dots, n\}$:

$$\hat{U}_{it} = \boldsymbol{\theta}'_i \hat{\mathbf{U}}_{-it} + V_{it}, \quad t \in \{1, \dots, T\}.$$

At the end of Steps 2 and 3, we can conduct the relevant inference on the structures of the covariance or partial covariance matrices. We can also provide updated prediction future outcomes. We detail those in the next subsection. Also note that the nonzero estimates of $\boldsymbol{\theta}_i$ shed light on the links among idiosyncratic components.⁵

2.3 Estimators and Inference Procedure

In a pure prediction exercise one is usually interested in the linear projection of Y_{it} onto $(\mathbf{X}'_{it}, \mathbf{F}'_t, \mathbf{U}'_{-it})'$, which results in the factor-augmented regression model (FARM)

$$Y_{it} = \gamma_i' \mathbf{X}_{it} + \boldsymbol{\lambda}_i' \mathbf{F}_t + \boldsymbol{\theta}_i' \mathbf{U}_{-it} + \varepsilon_{it}, \quad t \in \{1, \dots, T\}, \quad (2.2)$$

for each given i , and can be predicted by

$$\hat{Y}_{it} := \hat{\gamma}_i' \mathbf{X}_{it} + \hat{\boldsymbol{\lambda}}_i' \hat{\mathbf{F}}_t + \hat{\boldsymbol{\theta}}_i' \hat{\mathbf{U}}_{-it}; \quad i \in \{1, \dots, n\}. \quad (2.3)$$

This will be called `FarmPredict`. Note that model (2.2) is equivalent to using the predictors X_{it}, \mathbf{Y}_{-it} and \mathbf{F}_t , which augment predictors X_{it}, \mathbf{Y}_{-it} by using the common factors \mathbf{F}_t . The form

⁵The three-stage procedure describe here could be replaced by a single-step joint estimation. However, not only the computational burden will be much higher, but also the technical challenges will be greater. I believe that the simplicity of the method is more a blessing than a curse.

in (2.2) mitigates the collinearity issues in high dimensions.

Model (2.2) also bridges factor regression ($\boldsymbol{\theta}_i = 0$) on one end and (sparse) regression on the other end with $\boldsymbol{\lambda}_i = \boldsymbol{\Lambda}'_{-i}\boldsymbol{\theta}_i$, where $\boldsymbol{\Lambda}_{-i}$ is the loading matrix without the i^{th} row. In the latter case, model (2.2) becomes a (sparse) regression model:

$$Y_{it} = \boldsymbol{\gamma}_i' \mathbf{X}_{it} + \boldsymbol{\theta}_i' \mathbf{R}_{-it} + \varepsilon_{it}, \quad t \in \{1, \dots, T\}. \quad (2.4)$$

In this case, FARM specification as in (2.2) decorrelates the variables \mathbf{R}_{-i} in (2.4). It makes the model selection consistency much easier to satisfy and forms the basis of `FarmSelect` in (Fan et al., 2020). In general, for FARM (2.2) with sparsity, `FarmPredict` chooses additional idiosyncratic components to enhance the prediction of the factor regression.

In other applications, the structure of the idiosyncratic components $\mathbf{U} = (U_1, \dots, U_n)'$ is the objective of interest. An estimator for $\boldsymbol{\Sigma} = \mathbb{E}(\mathbf{U}_t \mathbf{U}_t')$ could be simply given by

$$\hat{\boldsymbol{\Sigma}} := \frac{1}{T} \sum_{t=1}^T \hat{\mathbf{U}}_t \hat{\mathbf{U}}_t'. \quad (2.5)$$

In order to properly understand the (linear) relation between a pair (U_{it}, U_{jt}) of \mathbf{U}_t , a simple covariance estimate sometimes is not enough. In applications, it is often desirable to have a direct measure of how U_{it} and U_{jt} are connected. By direct connection, we meant the relation between those units removing the contribution of other variables of \mathbf{U}_t . For this purpose, we use the partial covariance between U_{it} and U_{jt} , defined for any pair $i, j \in \{1, \dots, n\}$ as:

$$\pi_{ij} := \mathbb{E}(V_{ijt} V_{jit}),$$

where $V_{ijt} := U_{it} - \text{Proj}(U_{it} | \mathbf{U}_{-ij,t})$ and $\text{Proj}(U_{it} | \mathbf{U}_{-ij,t})$ denotes the linear projection of U_{it} onto the space spanned by all the units except i and j , which we denote by $\mathbf{U}_{-ij,t}$. We suggest to estimate the partial covariance matrix $\boldsymbol{\Pi} := (\pi_{ij})$ by

$$\hat{\boldsymbol{\Pi}} := (\hat{\pi}_{ij}) \quad \text{and} \quad \hat{\pi}_{ij} := \frac{1}{T} \sum_{t=1}^T \hat{V}_{ij,t} \hat{V}_{ji,t}, \quad (2.6)$$

where $\hat{V}_{ij,t}$ is the residual of the LASSO regression of \hat{U}_{it} onto $\hat{\mathbf{U}}_{-ij,t}$ for $i, j \in \{1, \dots, n\}$.

We also would like to conduct formal test on the population structure of \mathbf{U}_t . Specifically, we

propose a test for the following null hypothesis on the covariance matrix

$$\mathbb{H}_{\mathcal{D}}^{\Sigma} : \Sigma_{\mathcal{D}} = \Sigma_{\mathcal{D}}^0, \quad \mathcal{D} \subseteq \{1, \dots, n\} \times \{1, \dots, n\}, \quad (2.7)$$

for a given subset \mathcal{D} , where $\Sigma_{\mathcal{D}}$ denotes the elements of Σ indexed by \mathcal{D} and we allow $d := |\mathcal{D}|$ to diverge as $n, T \rightarrow \infty$. For example, to test if Σ is diagonal, \mathcal{D} consists of all off diagonal elements and $\Sigma_{\mathcal{D}}^0 = \mathbf{0}$. To test if Σ is block diagonal, \mathcal{D} can be taken to the corresponding off-diagonal blocks. Similarly, for testing the structure on the partial covariance matrix

$$\mathbb{H}_{\mathcal{D}}^{\Pi} : \Pi_{\mathcal{D}} = \Pi_{\mathcal{D}}^0, \quad \mathcal{D} \subseteq \{1, \dots, n\} \times \{1, \dots, n\}. \quad (2.8)$$

The null hypotheses (2.7) and (2.8) nest several cases of interest in applications. The most common would be to test for a diagonal or a block diagonal structure in Σ and/or Π . But it also accommodates other structures.⁶ The task of estimating Σ is well documented in literature even in high-dimensional setups; see, for example, Ledoit and Wolf (2004,2012,2017,2020), Fan et al. (2008), Lam and Fan (2009), or Fan et al. (2013).⁷

The challenges for testing (2.7) and (2.8) are similar and can be summarized as follows:

1. As we allow for both n and d to diverge to infinite as T grows, sometimes at a faster rate, we have a high-dimensional test where some sort of Gaussian approximation result for dependent data must be deployed as we also allow the number covariances to be tested (d) to diverge. In this case, a high-dimensional long-run covariance matrix must be estimated if one expects to get (asymptotic) correct test size.
2. We do not directly observe $\{\mathbf{U}_t\}$ or $\{V_{ij,t}\}$. Instead, we have an estimate of both from a postulated model on observable random variables. Therefore, the estimation error must be taken into account to claim some sort of asymptotic properties of the test. In fact, it is not uncommon to obtain estimates of both $\{\mathbf{U}_t\}$ and $\{V_{ij,t}\}$ from a multi-stage estimation procedure as we illustrate later in this paper.

We propose to test (2.7) using the statistic

$$S_{\mathcal{D}}^{\Sigma} := \|\sqrt{T}(\hat{\Sigma}_{\mathcal{D}} - \Sigma_{\mathcal{D}}^0)\|_{\max}. \quad (2.9)$$

⁶With minor changes, the proposed test can also be used to test the null $M\text{vec}(\Sigma) = \mathbf{m}$ for some $(d \times n^2)$ matrix M and d -dimensional vector \mathbf{m} where $d := d_T$ is also a function of T .

⁷See Ledoit and Wolf (2021a) for a recent survey or the book by Fan et al. (2020).

The quantiles of $S_{\mathcal{D}}^{\Sigma}$ are approximated by a Gaussian bootstrap. To describe the procedure, let $\mathbf{\Upsilon}_{\Sigma}$ denote the $(d \times d)$ covariance matrix for the vectorized submatrix $(\tilde{\sigma}_{ij})_{(i,j) \in \mathcal{D}}$, where $\tilde{\sigma}_{ij} := \frac{1}{T} \sum_{t=1}^T U_{i,t} U_{j,t}$. Since the process $\{\mathbf{U}_t\}$ might present some form of temporal dependence (refer to Assumption 3(c)) we estimate $\mathbf{\Upsilon}_{\Sigma}$ using a Newey-West-type estimator. For a given integrable function $k(\cdot)$ with $k(0) = 1$ and bandwidth $h > 0$, $\mathbf{\Upsilon}_{\Sigma}$ is estimated by

$$\hat{\mathbf{\Upsilon}}_{\Sigma} := \sum_{|\ell| < T} k(\ell/h) \hat{\mathbf{M}}_{\Sigma, \ell} \quad \text{and} \quad \hat{\mathbf{M}}_{\Sigma, \ell} := \frac{1}{T} \sum_{t=\ell+1}^T \hat{\mathbf{D}}_{\Sigma, t} \hat{\mathbf{D}}'_{\Sigma, t-\ell}, \quad (2.10)$$

where $\hat{\mathbf{D}}_{\Sigma, t}$ is a d -dimensional vector with entries given by $\hat{U}_{it} \hat{U}_{jt} - \hat{\sigma}_{ij}$ for $(i, j) \in \mathcal{D}$. Finally, let $c_{\Sigma}^*(\tau)$ be the τ -quantile of the Gaussian bootstrap

$$S_{\mathcal{D}}^* := \|\mathbf{Z}_{\Sigma}^*\|_{\infty}; \quad \mathbf{Z}_{\Sigma}^* | \mathbf{X}, \mathbf{Y} \sim \mathcal{N}(\mathbf{0}, \hat{\mathbf{\Upsilon}}_{\Sigma}).$$

Theorem 4 demonstrates the validity of Gaussian bootstrap procedure described above, i.e., it states conditions under which the τ -quantile of the test statistic (2.9) can be approximated by $c_{\Sigma}^*(\tau)$ in the appropriate sense.

Similarly, the test statistic for (2.8) is given by

$$S_{\mathcal{D}}^{\Pi} := \|\sqrt{T}(\hat{\mathbf{\Pi}}_{\mathcal{D}} - \mathbf{\Pi}_{\mathcal{D}}^0)\|_{\max}. \quad (2.11)$$

Let $\mathbf{\Upsilon}_{\Pi}$ denote the $(d \times d)$ covariance matrix of $(\tilde{\pi}_{ij})_{(i,j) \in \mathcal{D}}$ where $\tilde{\pi}_{ij} := \frac{1}{T} \sum_{t=1}^T V_{ij,t} V_{ji,t}$. For a given kernel $\mathcal{K}(\cdot) \in \mathbb{K}$ and bandwidth $h > 0$ where the class of kernels \mathbb{K} is described below in (3.10), $\mathbf{\Upsilon}_{\Pi}$ is estimated by

$$\hat{\mathbf{\Upsilon}}_{\Pi} := \sum_{|\ell| < T} \mathcal{K}(\ell/h) \hat{\mathbf{M}}_{\Pi, \ell}; \quad \hat{\mathbf{M}}_{\Pi, \ell} := \frac{1}{T} \sum_{t=\ell+1}^T \hat{\mathbf{D}}_{\Pi, t} \hat{\mathbf{D}}'_{\Pi, t-\ell}, \quad (2.12)$$

where $\hat{\mathbf{D}}_{\Pi, t}$ is a d -dimensional vector with entries given by $\hat{V}_{ij,t} \hat{V}_{ji,t} - \hat{\pi}_{ij}$ for $(i, j) \in \mathcal{D}$. Also, let $c_{\Pi}^*(\tau)$ be the τ -quantile of the Gaussian bootstrap

$$S_{\mathcal{D}}^* := \|\mathbf{Z}_{\Pi}^*\|_{\infty}; \quad \mathbf{Z}_{\Pi}^* | \mathbf{X}, \mathbf{Y} \sim \mathcal{N}(\mathbf{0}, \hat{\mathbf{\Upsilon}}_{\Pi}).$$

Theorem 5 demonstrate the validity of Gaussian bootstrap procedure describe above, i.e., it states conditions under which the τ -quantile of the test statistic (2.11) can be approximated by $c_{\Pi}^*(\tau)$ in the appropriate sense.

3 Theoretical Results

In this section we collect all the theoretical guarantees for the estimation of the model (2.1) by using the proposed three-stage method described above. Specifically, Section 3.1 deals with estimation and Section 3.2 with inference on the (partial) covariance structure of $\mathbf{\Pi}$.

To present the next results it is convenient to use a more compact notation. For each $i = 1, \dots, n$, we define the T -dimensional vectors $\mathbf{Y}_i := (Y_{i1}, \dots, Y_{iT})'$ and $\mathbf{U}_i := (U_{i1}, \dots, U_{iT})'$. We also define the $(T \times k)$ matrix of covariates $\mathbf{X}_i := (\mathbf{X}_{i1}, \dots, \mathbf{X}_{iT})'$ for each $i = 1, \dots, n$ and the $(T \times r)$ matrix of factors $\mathbf{F} := (\mathbf{F}_1, \dots, \mathbf{F}_T)'$ such that (2.1) can be represented as

$$\begin{aligned} \mathbf{Y}_i &= \mathbf{X}_i \boldsymbol{\gamma}_i + \mathbf{F} \boldsymbol{\lambda}_i + \mathbf{U}_i, & i = 1, 2, \dots, n, \\ &= \mathbf{X}_i \boldsymbol{\gamma}_i + \mathbf{R}_i, \end{aligned} \tag{3.1}$$

where $\mathbf{R}_i := \mathbf{F} \boldsymbol{\lambda}_i + \mathbf{U}_i$.

When no confusion is likely to arise, we also define for each $t = 1, \dots, T$, the n -dimensional vectors $\mathbf{Y}_t := (Y_{1t}, \dots, Y_{nt})'$ and $\mathbf{U}_t := (U_{1t}, \dots, U_{nt})'$; and the nk -dimensional vector $\mathbf{X}_t := (\mathbf{X}'_{1t}, \dots, \mathbf{X}'_{nt})'$. Also, set the $(n \times nk)$ block diagonal matrix $\mathbf{\Gamma}$ whose block diagonal is given by $(\boldsymbol{\gamma}'_1, \dots, \boldsymbol{\gamma}'_n)$ and the $(n \times r)$ loading matrix $\mathbf{\Lambda} := (\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_n)'$. Then, (2.1) can also be represented as

$$\begin{aligned} \mathbf{Y}_t &= \mathbf{\Gamma} \mathbf{X}_t + \mathbf{\Lambda} \mathbf{F}_t + \mathbf{U}_t, & t = 1, 2, \dots, T \\ &= \mathbf{\Gamma} \mathbf{X}_t + \mathbf{R}_t, \end{aligned} \tag{3.2}$$

where $\mathbf{R}_t := \mathbf{\Lambda} \mathbf{F}_t + \mathbf{U}_t$.

3.1 Estimation

For the factor model structure we consider the following set of assumptions

Assumption 2 (Factor Model). *Assume:*

- (a) $\mathbb{E}(\mathbf{F}_t) = \mathbf{0}$, $\mathbb{E}(\mathbf{F}_t \mathbf{F}_t') = \mathbf{I}_r$ and $\mathbf{\Lambda}' \mathbf{\Lambda}$ is a diagonal matrix;
- (b) All eigenvalues of $\mathbf{\Lambda}' \mathbf{\Lambda} / n$ are bounded away from zero and infinity as $n \rightarrow \infty$;
- (c) $\|\boldsymbol{\Sigma} - \mathbf{\Lambda} \mathbf{\Lambda}'\| = O(1)$; and
- (d) $\|\mathbf{\Lambda}\|_{\max} \leq C$.

Remark 2. *Assumption 2 is standard in the factor model literature. Note also that the assumption that $\mathbb{E}(\mathbf{F}_t) = \mathbf{0}$ is not restrictive as our approach considers a first-step estimation which may include a constant in the set of regressors. It is also needed for identifiability.*

In order to present the results in a unified manner for both light and heavy tail distributions, we state the next assumption in terms the Orlicz norm of the random variables. Specifically, since we are only concerned with polynomial and exponential tails we define the following subset of unbounded, convex, real-valued functions that vanish at the origin:

$$\begin{aligned} \Psi := \{ & \psi_p : \mathbb{R}^+ \rightarrow \mathbb{R}^+ : \psi_p(x) = x^p, p \geq 6 \\ & \text{or } \psi_p(x) = x\mathbf{1}[0 \leq x^p < (1-p)/p] + [\exp(x^p) - 1]\mathbf{1}[x^p \geq (1-p)/p], p > 0\}. \end{aligned} \quad (3.3)$$

Also, for each $\psi_p \in \Psi$, we define $\psi_{p+}(x) := x^{p+\epsilon}$ for some $\epsilon > 0$ if $\psi_p(x) = x^p$; otherwise (for the exponential case) $\psi_{p+} := \psi_p$.

Assumption 3 (Moments and Dependency). *There exists a constant $C < \infty$ and function $\psi_p \in \Psi$ defined in (3.3) such that, for all $i = 1, \dots, n$; $\ell = 1, \dots, k$; $s, t = 1, \dots, T$; and $j = 1, \dots, r$:*

$$(a) \quad \|X_{it\ell}\|_{\psi_{p+}} \leq C, \|U_{it}\|_{\psi_{p+}} \leq C, \|F_{j,t}\|_{\psi_{p+}} \leq C;$$

$$(b) \quad \|(\mathbf{X}'_i \mathbf{X}_i / T)^{-1}\|_{\max} \|_{\psi_{p+}} \leq C;$$

(c) *The process $\{(\mathbf{X}'_{S,t}, \mathbf{F}'_t, \mathbf{U}'_t)'\}, t \in \mathbb{Z}\}$ is weakly stationary with strong mixing coefficient α satisfying $\alpha(m) \leq \exp(-2cm)$ for some $c > 0$ and for all $m \in \mathbb{Z}$, where $\mathbf{X}_{S,t}$ denotes the vector \mathbf{X}_t after excluding all deterministic (non-random) components.*

$$(d) \quad \|n^{-1/2} [\mathbf{U}'_s \mathbf{U}_t - \mathbb{E}(\mathbf{U}'_s \mathbf{U}_t)]\|_{\psi_{p+}} \leq C;$$

$$(e) \quad \|n^{-1/2} \sum_{i=1}^n \lambda_{j,i} U_{it}\|_{\psi_{p+}} \leq C; \text{ and}$$

$$(f) \quad \log n = o\left(\frac{T^{p/4}}{[\log T]^2}\right).$$

A few words about Assumption 3 is in order. Assumptions (3.a) and (3.c) allow us to apply a Marcinkiewicz-Zygmund type inequality for partial sums to deal with the polynomial tails (Rio (1994) and Doukhan and Louhichi (1999)) and a Bernstein inequality (Merlevède et al. (2009) - Theorem 2) to control exponential tails. Moreover, Assumption (3.c) excludes the deterministic component of \mathbf{X}_t to accommodate possibly non-random non-stationary (but uniformly bounded by

(a)) covariates. Assumption (3.d) is only used to prove results for the first-stage estimation in case it is performed by ordinary least-squares (Theorem 1). Assumption (3.d) controls for the level of cross-sectional dependence among the units. As we allow the number of units to diverge with T , some sort of control on this quantity is necessary which is not implied by (3.c). Assumption (3.e) has a similar role to (3.d) but in terms of linear combinations of the the idiosyncratic components. Assumption (3.e) only bounds the growth rate of the number of units n to be sub-exponential with respect to T . As a matter of fact, this assumption is only binding in the exponential tail case, otherwise the rate conditions imposed in the theorems below imply (3.e).

For each $i = 1, \dots, n$, let $\mathbf{R}_i := \mathbf{F}\boldsymbol{\lambda}_i + \mathbf{U}_i$ denote the unobservable error term in (3.1), $\hat{\boldsymbol{\gamma}}_i$ the least-squares estimator of $\boldsymbol{\gamma}_i$ and $\hat{\mathbf{R}}_i := \mathbf{Y}_t - \mathbf{X}_t \hat{\boldsymbol{\gamma}}_i$ the vector of residuals. Also set $\hat{\mathbf{R}} := (\hat{\mathbf{R}}_1, \dots, \hat{\mathbf{R}}_n)'$ and $\mathbf{R} := (\mathbf{R}_1, \dots, \mathbf{R}_n)'$. We must control for the least-squares estimation error in the first step of the proposed methodology. The next result gives a bound for the maximum entry of the $(n \times T)$ matrix $\hat{\mathbf{R}} - \mathbf{R}$ when the first-stage is conducted by OLS in a linear setup.

Theorem 1. *Under Assumption 3(a)-(d)*

$$\max_{i,t} \|\hat{R}_{it} - R_{it}\|_{\psi_{p/4}} \leq \frac{C_{k,\psi}}{\sqrt{T}}$$

$$\|\hat{\mathbf{R}} - \mathbf{R}\|_{\max} = O_P \left[\frac{\psi_{p/4}^{-1}(nT)}{\sqrt{T}} \right],$$

where the $C_{k,\psi}$ is a constant only depending on k and ψ_p .

Remark 3. *In case the first step of the method involves more complicated estimation, we write $\|\hat{\mathbf{R}} - \mathbf{R}\|_{\max} = O_P(\omega)$, where $\omega := \omega_{n,T}$ is a non-negative sequence. This will be used in the next theorems.*

Define the $(n \times T)$ matrices $\mathbf{Y} := (\mathbf{Y}_1, \dots, \mathbf{Y}_T)$ and $\mathbf{U} := (\mathbf{U}_1, \dots, \mathbf{U}_T)$; and the $(nk \times T)$ matrix $\mathbf{X} := (\mathbf{X}_1, \dots, \mathbf{X}_T)$. We can write (2.1) in the matrix form as

$$\mathbf{Y} = \boldsymbol{\Gamma}\mathbf{X} + \boldsymbol{\Lambda}\mathbf{F}' + \mathbf{U}. \quad (3.4)$$

Notice that $\hat{\mathbf{R}} = \boldsymbol{\Lambda}\mathbf{F}' + \tilde{\mathbf{U}}$ where $\tilde{\mathbf{U}} := \mathbf{U} + \hat{\mathbf{R}} - \mathbf{R}$ and $(\boldsymbol{\Lambda}, \mathbf{F})$ can be estimated by Principal Component Analysis (PCA), which minimizes

$$q(\boldsymbol{\Lambda}, \mathbf{F}) := \|\hat{\mathbf{R}} - \boldsymbol{\Lambda}\mathbf{F}'\|_F^2, \quad (3.5)$$

with respect to $\mathbf{\Lambda}$ and \mathbf{F} , subject to the normalization $\mathbf{F}'\mathbf{F}/T = \mathbf{I}_r$. The solution $\hat{\mathbf{F}}$ is the matrix whose columns are \sqrt{T} times r eigenvectors of the top r eigenvalues of $\hat{\mathbf{R}}'\hat{\mathbf{R}}$ and $\hat{\mathbf{\Lambda}} = \hat{\mathbf{R}}\hat{\mathbf{F}}/T$.

Since we do not directly observe \mathbf{U} , in the third step of our estimation procedure we use $\hat{\mathbf{U}} := \hat{\mathbf{R}} - \hat{\mathbf{\Lambda}}\hat{\mathbf{F}}'$ instead. Therefore, we must control of the estimation error in the factor model given by $(n \times T)$ matrix $\hat{\mathbf{U}} - \mathbf{U}$ which is the main purpose of Theorem 2 below. Also, it is well know fact that the loading matrix $\mathbf{\Lambda}$ and the factors \mathbf{F} are not separably identified since $\mathbf{\Lambda}\mathbf{F}_t = \mathbf{\Lambda}\mathbf{H}'\mathbf{H}\mathbf{F}_t$ for any matrix \mathbf{H} such that $\mathbf{H}'\mathbf{H} = \mathbf{I}_r$. If we let $\mathbf{H} := T^{-1}\mathbf{V}^{-1}\hat{\mathbf{F}}'\mathbf{F}\mathbf{\Lambda}'\mathbf{\Lambda}$, where \mathbf{V} is the $(r \times r)$ diagonal matrix containing the r largest eigenvalues of $\hat{\mathbf{R}}'\hat{\mathbf{R}}/T$ in decreasing order, we have that $\mathbf{H}\mathbf{F}_t$ is identified as $\mathbf{\Lambda}\mathbf{F}_t$ is identified.

The result below first appeared in Bai (2003) for the case of fixed (n, T) , and was further extended to hold uniformly in $(i \leq n, t \leq T)$ by Fan et al. (2013). Fan et al. (2020) makes the conditions modular. However, both consider the case when the factor model is estimated using the true data as opposed to an “estimated” one as in our case. Therefore, the next result is a generalization that takes into account that pre-estimation error term.

Theorem 2. *Let $\omega := \omega_{n,T}$ be a non-negative sequence such that $\|\hat{\mathbf{R}} - \mathbf{R}\|_{\max} = O_P(\omega)$. Then, under Assumptions 1 -3 and $\psi_p^{-1}(n^2)/\sqrt{T} + \psi_p^{-1}(nT)\omega = O(1)$, we have that*

(a)

$$\max_{t \leq T} \|\hat{\mathbf{F}}_t - \mathbf{H}\mathbf{F}_t\|_2 = O_P \left[\frac{1}{\sqrt{T}} + \frac{\psi_p^{-1}(T)}{\sqrt{n}} + \omega\psi_{p/2}^{-1}(nT) \right],$$

(b)

$$\max_{i \leq n} \|\hat{\boldsymbol{\lambda}}_i - \mathbf{H}\boldsymbol{\lambda}_i\|_2 = O_P \left[\frac{\psi_{p/2}^{-1}(n)}{\sqrt{T}} + \frac{1}{\sqrt{n}} + \omega \right],$$

(c)

$$\|\hat{\mathbf{U}} - \mathbf{U}\|_{\max} = O_P \left[\frac{\psi_p^{-1}(n)\psi_p^{-1}(T)}{\sqrt{T}} + \frac{\psi_p^{-1}(T)}{\sqrt{n}} + \omega\psi_{p/2}^{-1}(nT) \right].$$

By setting $\omega = 0$, i.e., no estimation error in the first step, we recover Theorem 4 and Corollary 1 in Fan et al. (2013) under sub-Gaussian assumption. Also it is important to notice that in order to have the error $\|\hat{\mathbf{U}} - \mathbf{U}\|_{\max}$ vanishing in probability we must have the pre-estimation error $\|\hat{\mathbf{R}} - \mathbf{R}\|_{\max}$ of order (in probability) smaller than $1/\psi_{p/2}^{-1}(nT)$.

We have decided not to replace ω in Theorem 2 with the rate obtained in Theorem 1 as the latter only applies to the least square estimator. In some applications, however, the first step of the procedure could be done using a different type of estimator. For instance a penalized adaptive

Huber regression (Fan et al., 2017) if the number of features k is comparable or even larger than T and the tail of the distribution is heavy. By stating the Theorem 2 in terms of a generic rate, it is easier to account for the effect of a different estimator. By combining Theorem 1 and 2 we have the following corollary

Corollary 1. *Under the same assumptions of Theorems 1 and 2, for the OLS used in the first-stage to obtain $\hat{\mathbf{R}}$, we have*

$$\|\hat{\mathbf{U}} - \mathbf{U}\|_{\max} = O_P \left[\frac{\psi_p^{-1/6}(nT)}{\sqrt{T}} + \frac{\psi_p^{-1}(T)}{\sqrt{n}} \right].$$

In particular for the sub-Gaussian case ($\psi(x) = \exp(x^2) - 1$) we have

$$\|\hat{\mathbf{U}} - \mathbf{U}\|_{\max} = O_P \left[\frac{[\log(nT)]^3}{\sqrt{T}} + \sqrt{\frac{\log T}{n}} \right],$$

and for polynomial tails ($\psi(x) = x^p$)

$$\|\hat{\mathbf{U}} - \mathbf{U}\|_{\max} = O_P \left[\frac{n^{6/p}}{T^{1/2-6/p}} + \frac{T^{1/p}}{\sqrt{n}} \right].$$

For notational convenience, for each $i \in \{1, \dots, n\}$, consider the split $\mathbf{U}' = (\mathbf{U}_i, \mathbf{U}_{-i})$ where \mathbf{U}_i is a T -dimensional vector and \mathbf{U}_{-i} a $T \times (n-1)$ -dimensional matrix. Analogously, we split $\hat{\mathbf{U}}' = (\hat{\mathbf{U}}_i, \hat{\mathbf{U}}_{-i})$. Then for a the penalized parameter $\xi \geq 0$, the LASSO objective function can be written for each $i \in \{1, \dots, n\}$

$$L(\boldsymbol{\theta}) + \text{Penalty}(\boldsymbol{\theta}) := \frac{1}{T} \|\hat{\mathbf{U}}_i - \hat{\mathbf{U}}_{-i} \boldsymbol{\theta}\|^2 + \xi \|\boldsymbol{\theta}\|_1. \quad (3.6)$$

To ensure a consistent estimation of $\boldsymbol{\theta}$, a sort of restricted strong convexity of the objective function is required when $n > T$. This in turns is ensured, in the case of a quadratic loss, by bounding the minimum eigenvalue on $\hat{\mathbf{U}}_{-i}' \hat{\mathbf{U}}_{-i} / T$ away from zero restrict to a cone (refer to Negahban et al. (2012) or Fan et al. (2020) for a thorough discussion). Here, we adopt the compatibility constant defined in van de Geer and Bühlmann (2009). For an index $\mathcal{S} \subseteq \{1, \dots, n\}$ and any n -dimensional vector \mathbf{v} , let $\mathbf{v}_{\mathcal{S}}$ be the vector containing only the elements of the vector \mathbf{v} indexed by \mathcal{S} . Thus, $\#\mathbf{v}_{\mathcal{S}} = \#\mathcal{S}$ and $\mathcal{S}^c := \mathcal{S} \setminus \{1, \dots, n\}$ is the complement of \mathcal{S} .

Definition 1. *For an $n \times n$ matrix \mathbf{M} , a set $\mathcal{S} \subseteq \{1, \dots, n\}$ and a scalar $\zeta \geq 0$, the compatibility*

constant is given by

$$\kappa(\mathbf{M}, \mathcal{S}, \zeta) := \inf \left\{ \frac{\|\mathbf{x}\|_{\mathbf{M}} \sqrt{|\mathcal{S}|}}{\|\mathbf{x}_{\mathcal{S}}\|_1} : \mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}_{\mathcal{S}^c}\|_1 \leq \xi \|\mathbf{x}_{\mathcal{S}}\|_1 \right\}, \quad (3.7)$$

where $\|\mathbf{x}\|_{\mathbf{M}} = \mathbf{x}'\mathbf{M}\mathbf{x}$. Moreover, we say that $(\mathbf{M}, \mathcal{S}, \zeta)$ satisfies the compatibility condition if $\kappa(\mathbf{M}, \mathcal{S}, \zeta) > 0$.

Notice that the square of the compatibility constant is close related to the minimum of the ℓ_1 -norm of the eigenvalues of Σ restricted to a cone in \mathbb{R}^n .

Theorem 3. Let $\eta := \eta_{n,T}$ be a non-negative sequence such that $\|\hat{\mathbf{U}} - \mathbf{U}\|_{\max} = O_P(\eta)$ and consider Assumption 3. For every $\epsilon > 0$ there is a constant $0 < C < \infty$ such that if the penalty parameter is set to

$$\xi = C \left[\frac{\psi_{p/2}^{-1}(n)}{\sqrt{T}} + \eta \psi_p^{-1}(T) \right] \quad (3.8)$$

and $s_0 := \max_{i \leq n} |\mathcal{S}_{0,i}|$ where $\mathcal{S}_{0,i} := \{j : \theta_{i,j} \neq 0\}$ obeys

$$s_0 = O \left[\kappa_0 \left(\eta [\psi_p^{-1}(nT) + \eta] + \frac{\psi_{p/2}^{-1}(n^2)}{\sqrt{T}} \right)^{-1} \right], \quad (3.9)$$

with $\kappa_0 := \min_{i \leq n} \kappa_i$ and $\kappa_i := \kappa[\mathbb{E}(\mathbf{U}'_{-i}\mathbf{U}_{-i})/T, \mathcal{S}_{0,i}, 3]$ defined in (3.7). Then, for any minimizer $\hat{\boldsymbol{\theta}}_i$ of (3.6), with probability at least $1 - \epsilon$:

$$T^{-1}(\hat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i)' \mathbf{U}'_{-i} \mathbf{U}_{-i} (\hat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i) + \xi \|\hat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i\|_1 \leq 8 \frac{\xi^2 s_0}{\kappa_0}. \quad i \in \{1, \dots, n\},$$

where the left right side is taken to be $+\infty$ whenever $\kappa_0 = 0$.

Remark 4. Notice that we apply the compatibility condition on the non-random covariance matrix $\mathbb{E}(\mathbf{U}'_{-i}\mathbf{U}_{-i})/T$ instead of the estimated random covariance matrix $\hat{\mathbf{U}}'_{-i}\hat{\mathbf{U}}_{-i}/T$ or the “unobservable” random matrix $\mathbf{U}'_{-i}\mathbf{U}_{-i}/T$. Careful review of the proofs reveals that the same is true for the gradient of the objective function that defines our parameter via a first order condition.

Once again, we purposely avoided to replace η in Theorem 3 with the rate of Corollary 1 to make it readily applicable to the case when a different type of factor models was used or, as a matter of fact, any other pre-estimation procedure. By plugging the rate of Corollary 1 into η we have the next corollary

Corollary 2. *If η defined in Theorem 3 is taken to be rate given by Corollary 1 and the compatibility condition holds, i.e.: $\kappa_0 \geq C > 0$ then under the conditions of the Theorem 3:*

$$\max_{i \leq n} \|\hat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i\|_1 = O_P \left[\left(\frac{\psi_p^{-1}(T)\psi_{p/6}^{-1}(nT)}{\sqrt{T}} + \frac{\psi_{p/2}^{-1}(T)}{\sqrt{n}} \right) s_0 \right].$$

3.2 Inference

We now obtain the null distributions of our test statistics for the structures of the covariance and the partial covariance. Recall the setup and notation of section 2.3. In particular, we consider the kernel $k(\cdot)$ appearing in the covariance estimator defined by (2.10) belongs to the class defined in Andrews (1991) which we reproduce below for convenience

$$\mathbb{K} := \{f : \mathbb{R} \rightarrow [-1, 1] : f(0) = 1, f(x) = f(-x), \forall x \in \mathbb{R}, \int f^2(x)dx < \infty, f \text{ is continuous}\}. \quad (3.10)$$

It includes most of the well-known kernel used in density estimation literature such as the truncated, Bartlett, Parzen, Quadratic Spectral, Tukey-Hanning among others. To avoid confusion, it is worth to point out that our tuning parameter h , also called bandwidth parameter by Andrews (1991), is supposed to diverge, as opposed to the bandwidth in the density kernel estimation setup, which is expected to shrink towards zero.

Theorem 4. *Let $\eta := \eta_{n,T}$ and $\nu := \nu_{n,T}$ be non-negative sequence such that $\|\hat{\mathbf{U}} - \mathbf{U}\|_{\max} = O_P(\eta)$ and $\max_{i,t} \|\hat{R}_{it} - R_{it}\|_{\psi_p} = O(\nu)$ and $\mathcal{K} \in \mathbb{K}$. Under Assumptions 1-3, if further*

(a) $\{\mathbf{U}_t\}$ is fourth-order stationary process

(b) $\|\text{diag}(\boldsymbol{\Upsilon}_\Sigma)\|_\infty \geq \underline{c}$ for some $\underline{c} > 0$

(c) As $h, n, T \rightarrow \infty$:

$$(c.1) \quad \frac{(\log n)^{7/6}\psi_{p/2}^{-1}(n)}{T^{1/6}} + \frac{\sqrt{\log T \log n}\psi_{p/2}^{-1}(n)\psi_{p/2}^{-1}(T^{1/4})}{T^{1/4}} = o(1)$$

$$(c.2) \quad (\log n)^3 h \left[\eta(\psi_p^{-1}(nT))^3 + \psi_{p/4}^{-1}(n^4)/\sqrt{T} \right] = o(1)$$

$$(c.3) \quad (\log n)^3 \left(\sqrt{T}\eta^2 + \frac{r_1}{\sqrt{T}} + \frac{r_2}{\sqrt{n}} + r_3\nu \right) = o(1),$$

where the rates r_1, r_2, r_3 are defined in Lemma B.10 and $h > 0$ is the bandwidth parameter of the

covariance estimator defined in (2.10); then

$$\|\hat{\Upsilon}_\Sigma - \Upsilon_\Sigma\|_{\max} = O_P \left[h \left(\eta [\psi_p^{-1}(nT)]^3 + \psi_{p/4}^{-1}(n^4)/\sqrt{T} \right) \right] = o(1),$$

and

$$\sup_{\mathcal{D}} \sup_{\tau \in (0,1)} |\mathbb{P}[S_{\mathcal{D}}^\Sigma \leq c_\Sigma^*(\tau)] - \tau| = o(1),$$

where the first supremum is over all null hypotheses of the form (2.7) indexed by $\mathcal{D} \in \{1 : n\} \times \{1 : n\}$.

Remark 5. The rate assumptions (c.1)-(c.3) in Theorem 4 seem over complicated. However, they are a direct consequence of having the first and second step estimation error rates, ν and η respectively, explicitly appearing in the final rate and the general tail condition through the $\psi_p(\cdot)$ function. It allows the practitioner to directly adjust the final rate should (s)he prefer to employ different intermediate estimators. For instance, a LASSO estimator in the first step in case the number of covariates k is large enough or estimate the factor model by PCA variants. If we were to specialized to the sub-Gaussian case and incorporate the rates obtain in Theorem 1 and Corollary 1 we have the following Corollary

Corollary 3. Consider the sub-Gaussian where $\psi_2(x) = \exp(x^2)$. Suppose that the Assumptions 1-3 and conditions (a) and (b) of Theorem 4 hold. If the rates ν and η are set to be rates given by Theorem 1 and Corollary 1, respectively, then the conclusion of Theorem 4 holds provided that as $h, n, T \rightarrow \infty$:

$$(a) \log n = o(T^{1/18})$$

$$(b) h \left[\frac{(\log n)^{15/2}}{\sqrt{T}} + \frac{(\log n)^5}{\sqrt{n}} \right] = o(1)$$

$$(c) \frac{(\log n)^3 (\log T) \sqrt{T}}{n} = o(1).$$

Remark 6. Careful review of its proof reveals that (d.1) traces back to the Gaussian Approximation of the (unobservable) process $\left\{ \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{U}_t \mathbf{U}_t' - \mathbb{E} \mathbf{U}_t \mathbf{U}_t' \right\}_{T \geq 1}$; whereas (d.3) controls for the difference between $\mathbf{U}_t - \hat{\mathbf{U}}_t$ and, therefore, takes into account the estimation error of the first and second steps. Note the presence of ν and η in (d.3) which are absent in (d.1) Finally, (d.2) make sure that the bootstrap constructed in terms of the estimated covariance matrix is close to the bootstrap based in the true covariance. Note the presence of the bandwidth parameter h in (d.2).

Remark 7. In order to establish the rate of convergence in the last result of Theorem 4 we need an upper bound on the tails of the pre-estimation error namely $\|\widehat{\mathbf{Z}} - \mathbf{Z}\|_{\max}$. In fact, we need to control the tails of the factor model estimation to establish uniform bounds on $\|\widehat{U}_{it} - U_{it}\|_{\psi}$, which translate into obtain bounds on $\max_{jt} \|\widehat{F}_{jt} - F_{jt}\|_{\psi}$ and $\max_{ji} \|\widehat{\lambda}_{ji} - \lambda_{ji}\|_{\psi}$.

Theorem 5. Let $\eta := \eta_{n,T}$ and $\nu := \nu_{n,T}$ be non-negative sequence such that $\|\widehat{\mathbf{U}} - \mathbf{U}\|_{\max} = O_P(\eta)$ and $\max_{i,t} \|\widehat{R}_{it} - R_{it}\|_{\psi_p} = O(\nu)$ and $\mathcal{K} \in \mathbb{K}$ defined by (3.10). Under Assumptions 2-4 and the LASSO regularization parameter as in 3.8, if further

(a) $\{\mathbf{U}_t\}$ is fourth-order stationary process

(b) $\|\text{diag}(\mathbf{\Upsilon}_{\Pi})\|_{\infty} \geq \underline{c}$ for some $\underline{c} > 0$

(c) As $n, T \rightarrow \infty$:

$$(c.1) \quad \frac{(\log n)^{7/6} \psi_{p/2}^{-1}(n)}{T^{1/6}} + \frac{\sqrt{\log T \log n} \psi_{p/2}^{-1}(n) \psi_{p/2}^{-1}(T^{1/4})}{T^{1/4}} = o(1)$$

$$(c.2) \quad (\log n)^3 h \left(s_0 [\eta + \xi \psi_p^{-1}(n)] [s_0 \psi_p^{-1}(nT)]^3 + s_0 \frac{\psi_{p/4}^{-1}(n^4)}{\sqrt{T}} \right) = o(1)$$

$$(c.3) \quad (\log n)^3 \left(s_0^2 \left\{ \frac{r_1}{\sqrt{T}} + \frac{r_2}{\sqrt{n}} + r_3 \nu + \xi \psi_p^{-1}(n) + \sqrt{T} [\eta + \xi \psi_p^{-1}(n)]^2 \right\} \right) = o(1),$$

where the rates r_1, r_2, r_3 are defined in Lemma B.10, $\mathcal{K}(\cdot)$ and $h > 0$ is the bandwidth parameter of the covariance estimator defined in (2.12); then

$$\|\widehat{\mathbf{\Upsilon}}_{\Pi} - \mathbf{\Upsilon}_{\Pi}\|_{\max} = O_P \left(h \left\{ s_0 [\eta + \xi \psi_p^{-1}(n)] [s_0 \psi_p^{-1}(nT)]^3 + s_0 \frac{\psi_{p/4}^{-1}(n^4)}{\sqrt{T}} \right\} \right) = o(1)$$

and

$$\sup_{\mathcal{D}} \sup_{\tau \in (0,1)} |\mathbb{P}(S_{\mathcal{D}}^{\Pi} \leq c_{\Pi}^*(\tau)) - \tau| = o(1) \quad \text{under } \mathbb{H}_0^{\Pi},$$

where the first supremum is over all null hypotheses of the form (2.8) indexed by $\mathcal{D} \in \{1 \times n\} \times \{1 \times n\}$.

Remarks and Corollary analogous to Remarks 5-7 and Corollary 3 after Theorem 4 apply to Theorem. 5.

Remark 8. As opposed to the case of testing covariance, when testing partial covariance in high-dimensional setup, the sparse structure plays a role in terms of s_0 appearing in the rates (d.2) and (d.3). Therefore, these assumptions restricts the cases when the proposed partial covariance test has the correct asymptotic size. For instance, in the case of a complete dense partial covariance structure, i.e, all the regressors are active in all LASSO regressions we are likely to have s_0 of order of n and, therefore, (d.2) and (d.3) are not expected to hold.

4 Guide to Practice

As described before the methodology in this paper involves three steps. The first step consists of identifying known covariates that we may want to control for. This first step may involve the removal of deterministic trends and seasonal effects, for instance. This can be done either by parametric or nonparametric regressions. It is important to notice, however, that the convergence rates of the estimations in the subsequent steps will be influenced by the convergence rate of the estimation in the first part of the procedure.

After the data is filtered in the first step, one can test for remaining covariance structure. For instance, if the covariance matrix of the filtered data is (almost) diagonal, there is no need to estimate a latent factor structure and the practitioner may jump directly to the third step of the method.

On the other hand, if the covariance of the first-step filtered data is dense, a latent factor model should be considered and the number of factors must be determined. There are a number of methods proposed in the literature to achieve this goal. In this paper we consider either the eigenvalue ratio test of Horenstein (2013) or the information criteria put forward in Bai and Ng (2002). The factors can be estimated by the usual methods.

The last step involves a sparse regression in order to estimate any remaining links between idiosyncratic components. Before running the last step, the practitioner may test for a diagonal covariance matrix of the idiosyncratic terms. If the null is not rejected, there is no need for additional estimation. In case of rejection, the user can proceed with a LASSO regression. We recommend that the penalty term of the LASSO is selected by Bayesian Information Criterion (BIC) as advocated by Medeiros and Mendes (2016).

Finally, we would like to include a remark about the estimation of the long-run matrices when constructing the statistics for the tests of no remaining covariance structure. Usual methods discussed in the literature can be used here to select the kernel and the bandwidth. In the paper we use the simple Bartlett kernel with bandwidth given as $\lfloor T/3 \rfloor$.

5 Simulation

In this section we report simulation results to assess the finite-sample performance of the methodology depicted in this paper. The simulations are divided into two parts. In the first one, we evaluate

the finite-sample properties of the test for remaining covariance structure. In the second part, we highlight the informational gains when considering both the common factors and the idiosyncratic component.

We simulate 1,000 replications of the following model for various combinations of sample size (T) and number of variables (n):

$$\mathbf{Y}_{it} = \boldsymbol{\Lambda}'_i \mathbf{F}_t + W_{it}, \quad (5.1)$$

$$\mathbf{F}_t = 0.8 \mathbf{I}_r + \mathbf{E}_t, \quad (5.2)$$

$$W_{it} = \phi W_{it-1} + U_{it}, \quad (5.3)$$

$$U_{it} = \begin{cases} \theta_{12}U_{2t} + \theta_{13}U_{3t} + \theta_{14}U_{4t} + \theta_{15}U_{5t} + O_{it} & \text{if } i = 1 \\ O_{it} & \text{otherwise,} \end{cases} \quad (5.4)$$

where $\{O_{it}\}$ is a sequence of independent Gaussian random variables with zero mean and variance equal to 0.25, $\{\mathbf{E}_t\}$ is a sequence of r -dimensional independent random vectors normally distributed with zero mean and identity covariance, and \mathbf{I}_r is an $r \times r$ identity matrix. Furthermore, $\{O_{it}\}$ and $\{\mathbf{E}_t\}$ are mutually independent for all time periods, factors and variables. For each Monte Carlo replication, the vector of loadings is sampled from a Gaussian distribution with mean -6 and standard deviation 0.2 for $i = 1$ and mean 2 and unit variance for $i = 2, \dots, n$. The value of ϕ is either 0 or 0.5. The coefficients θ_{12} , θ_{13} , θ_{14} , and θ_{15} are equal to zero or 0.8, 0.9, -0.7, and 0.5, respectively. We set the true number of factor to be $r = 3$.

5.1 Test for Remaining Covariance Structure

We start by reporting results for the test of no remaining structure on the covariance matrix of $\mathbf{U}_t = (U_{1t}, \dots, U_{nt})'$. The null hypothesis considered is that all the covariances between the first variable ($i = 1$) and the remaining ones are all zero. For size simulations we set $\theta_{12} = \theta_{13} = \theta_{14} = \theta_{15} = 0$ in the DGP. In order to evaluate the effects of factor estimation as well as the methods in selecting the number of factors, we consider the following scenarios: (1) factors are known and there is no estimation involved; (2) factors are estimated by principal components but the number of factors are known; (3) the number of factors is determined by the eigenvalue ratio procedure of Horenstein (2013); (4)-(7) the number of factors is determined by one of the four information criteria proposed

by Bai and Ng (2002) as defined as

$$\begin{aligned} \text{IC}_1 &= \log[S(r)] + r \frac{n+T}{nT} \log\left(\frac{nT}{n+T}\right) \\ \text{IC}_2 &= \log[S(r)] + r \frac{n+T}{nT} \log C_{nT}^2 \\ \text{IC}_3 &= \log[S(r)] + r \frac{\log C_{nT}^2}{C_{nT}^2} \\ \text{IC}_4 &= \log[S(r)] + r \frac{(n+T-k) \log(nT)}{nT}. \end{aligned}$$

where $S(r) = \frac{1}{nT} \|\mathbf{R} - \hat{\mathbf{\Lambda}}_r \hat{\mathbf{F}}_r\|_2^2$ and $C_{nT} := \sqrt{\min(n, T)}$.

Tables 1 and 2 report the results of the empirical size of test for different significance levels. We consider the case of $\phi = 0$ in Table 1 and $\phi = 0.5$ in Table 2. The tables present the results when the factors are known in panel (a), the factors are unknown but the number of factors is known in panel (b), or the number of factors are estimated either by the information criterion IC_1 in panel (c) or the eigenvalue ratio procedure in panel (d). Table ?? in the Supplementary Material shows the results of the test when the number of factors are determined by $\text{IC}_2 - \text{IC}_4$.

A number of facts emerge from the inspection of the results in the Table 1. First, size distortions are small when the factors are known. In this case, the test is undersized when the pair (n, T) is small. When the factor are not known but the true number of factors is available, the size distortions are high only when $T = 100$ and $n = 50$ due to inaccurate estimation of factors. However, the distortions disappear when the pair (T, n) grows. In this case, the empirical size is similar to the situation reported in Panel (a). The finite performance of test in the case where the number of factors is selected by information criterion IC_1 is almost indistinguishable to the case reported in Panel (b). However, the results with the eigenvalue ratio procedure are much worse when $T = 100$ and $n = 50$. In this case, the procedure selects less factors than true number $r = 3$. For instance, the procedure selects 2 or less factors in 36% of the replications. Just as comparison, for $T = 100$ and $n = 50$, IC_1 underdetermines the number of factors only in 3.10% of the cases. For all the other combinations of T and n all the data-driven methods selects the correct number of factors in almost all replications.

When the idiosyncratic components are autocorrelated the size distortions are higher, as reported in Table 2. This is mainly caused by the well-known difficulties in the estimation of the long-run covariance matrix.

Tables 3–4 report the results of the empirical power. For evaluate power properties we set

$\beta_{12} = 0.8$, $\beta_{13} = 0.9$, $\beta_{14} = -0.7$, and $\beta_{15} = -0.5$ in the DGP. When the factors are known, the test always rejects the null and the empirical power is one for any significance level. On the other hand, when factors must be estimated but the number of factors are known, the power decreases as depicted in panel (b) in the tables. Nevertheless for $T = 500, 700$ the power is reasonably high, specially when test is conducted at a 10% significance level. For $T = 100$, the performance deteriorates as n grows. The results are similar when data-driven procedures are used to determine the number of factors. Finally, the conclusions are mostly the same whenever $\phi = 0$ or $\phi = 0.5$.

The main message of the simulation exercise is that the finite-sample performance of the proposed tests depend on the correct selection of factors. Nevertheless, for the DGP considered here, the usual data-driven methods available in the literature to determine the true number of factors seem to work reasonably well.

5.2 Informational Gains

The goal of this simulation is to compare, in a prediction environment, the three-stage method developed in the paper by evaluating the information gains in predicting Y_{1t} by three different methods. First, the predictions are computed from a LASSO regression of Y_{1t} on all the other $n - 1$ variables. This is the Sparse Regression (SR) approach. Second, we consider a principal component regression (PCR), i.e., an ordinary least squares (OLS) regression of the variable of interest on factors computed from the pool of other variables. Finally, we consider predictions constructed from the method proposed here, the `FarmPredict` methodology. Table 5 presents the results. The table presents the average mean squared error (MSE) over 5-fold cross-validation (CV) subsamples. As in the size and power simulations, we consider different combinations of T and n . We report results for the case where $\theta_{12} = 0.8$, $\theta_{13} = 0.9$, $\theta_{14} = -0.7$, and $\theta_{15} = -0.5$ in the DGP.

According to the DGP, the theoretical MSE is 0.25 when all the information is used. When just a factor is used, the MSE is 2.21. From the table is clear that there are significant informational gains when we consider both factors and the cross-dependence between idiosyncratic components. Several conclusions emerge from the table. First, it is clear that when the sample size increases the MSE reduces. This is expected. Second, the PCR's MSE is close to the theoretical value when the sample increases. The performance of the `FarmPredict` is quite remarkable when $T = 500$ or $T = 700$ and is always superior to Sparse Regression.

6 Applications

In this section we consider two applications with actual data to illustrate the benefits of the methodology developed in the paper. The first application deals with factor structure of asset returns, whereas the second one is about macroeconomic forecasting in data-rich environments.

6.1 Factor Models Network Structure in Asset Returns

6.1.1 The Dataset

We illustrate the methodology developed in this paper by studying the factor structure of asset returns. We consider monthly close-to-close excess returns from a cross-section of 9,456 firms traded in the New York Stock Exchange. The data starts on November 1991 and runs until December 2018. There are 326 monthly observations in total. In addition to the returns we also consider 16 monthly factors: Market (MKT), Small-minus-Big (SMB), High-minus-Low (HML), Conservative-minus-Aggressive (CMA), Robust-minus-Weak (RMW), earning/price ratio, cash-flow/price ratio, dividend/price ratio, accruals, market beta, net share issues, daily variance, daily idiosyncratic variance, 1-month momentum, and 36-month momentum. The firms are grouped according to 20 industry sectors as in Moskowitz and Grinblatt (1999). The following sectors are considered:⁸ Mining (602), Food (208), Apparel (161), Paper (81), Chemical (513), Petroleum (48), Construction (68), Primary Metals (133), Fabricated Metals (186), Machinery (710), Electrical Equipment (782), Transportation Equipment (166), Manufacturing (690), Railroads (25), Other transportation (157), Utilities (411), Department Stores (67), Retail (1018), Financial (3419), and Other (11).

6.1.2 Results

We start the analysis by looking at the correlation matrix of a sample of nine different sectors, namely: Mining, Food, Petroleum, Construction, Manufacturing, Utilities, Department Stores, Retail, and Financial. Figure 1 plots the correlations that are larger than 0.15 in absolute value. We also test for the null of diagonal covariance matrix. The null hypothesis is strongly rejected with p -value much lower than 1%. To conduct the test of the covariance matrix we use the simple sample estimator as described in the paper. However, the correlations plotted in Figure 1 and in the subsequent ones are based on the nonlinear shrinkage estimator proposed by Ledoit and Wolf (2020).

⁸The number between parenthesis indicate the number of firms in our sample that belong to each sector.

We proceed by regressing the daily returns on the observed 16 factors. These three factors explain most of the variation of the returns. Figure 2 shows the empirical distribution of the OLS estimates of factor loadings over the 9,456 regressions. Figure 3 presents the estimated correlations for the first-stage residuals. We focus on the nine sectors as before. The first-stage regression is efficient in removing the correlation within specific sectors in some cases. The most notable ones are Financial and Retail, followed by Construction, Petroleum, and Manufacturing. Nevertheless, the tests for diagonal covariance matrix reject the null even in these specific cases.

The second step is to conduct a principal component analysis on the residuals of the first-stage. The eigenvalue ratio procedure selects two factors, while all four information criteria points to a single factor. We proceed with two factors. Note that, by construction, the principal component factors are orthogonal to all the 16 risk factors considered in the first stage. Figure 4 shows the estimated correlations for the residuals of the second-stage. The latent factors are not able to reduce the correlations within each sector. However, when we consider the partial correlations the conclusions are much different. As can be seen from Figure 5 that the partial correlation matrices are (almost) diagonal. In addition, we are not able to reject the null of a diagonal covariance matrix at a 5% significance level.

Finally, in order to shed some light on the links among different sectors, we report how often that variables from sector i are selected in the third-stage LASSO regression for firms in sector j . The numbers are normalized by the total number of firms in each sector and are presented in Figure 6. The most interesting fact is that covariates from the financial sector are the ones most frequently selected for all the other sectors. This may indicate that there is a “financial factor” that was unmodeled in the first two stages.

The results presented here can be useful in applications where forecasting future returns is the goal, for instance. The results indicate that the inclusion of the returns of firms belonging to the financial sector may improve the performance of forecasting models. For example, if we run a regression of the residuals of the first-stage regression of firms that do not belong to the financial sector on the first principal component computed with the first-stage residuals only from the financial sector, we find a statistically significant coefficient in 28% of the cases.

6.2 Macroeconomic Forecasting

The second application consists of forecasting of a large set of monthly macroeconomic variables. We compare four different models: (1) Autoregressive model; (2) Sparse LASSO Regression (SR); (3) Principal Component Regression (PCR); and (4) a method based on the results in this paper (`farmPredict`).

6.2.1 The Dataset

Our data consist of variables from the FRED-MD database, which is a large monthly macroeconomic dataset designed for empirical analysis in data-rich macroeconomic environments. The dataset is updated in real time through the FRED database and is available from Michael McCracken's webpage.⁹ For further details, we refer to McCracken and Ng (2016).

We use the vintage as of October 2020. Our sample extends from January 1960 to December 2019 (719 monthly observations), and only variables with all observations in the sample period are used (119 variables). The dataset is divided into eight groups: (i) output and income; (ii) labor market; (iii) housing; (iv) consumption, orders and inventories; (v) money and credit; (vi) interest and exchange rates; (vii) prices; and (viii) stock market. Finally, all series are transformed in order to become approximately stationary as in McCracken and Ng (2016).

6.2.2 Setup and Methodology

In order to highlight the gains of exploring all relevant information in the the dataset, we construct one-step forecasts for each one of the 119 variables in the dataset according to the following models:

1. Autoregressive model (AR):

$$\hat{Y}_{i,t+1|t}^{(\text{AR})} = \hat{\phi}_{i0} + \hat{\phi}_{i1}\hat{Y}_{i,t} + \dots + \hat{\phi}_{ip}\hat{Y}_{i,t-p+1}, \quad i = 1, \dots, n,$$

where $\hat{\phi}_{i0}, \hat{\phi}_{i1}, \dots, \hat{\phi}_{ip}$, $i = 1, \dots, n$, are OLS estimates. This model will be also the first-stage model in our methodology.

2. AR + Sparse regression (SR):

$$\hat{Y}_{i,t+1|t}^{(\text{SR})} = \hat{Y}_{i,t+1|t}^{(\text{AR})} + \hat{R}_{i,t+1|t},$$

⁹<https://research.stlouisfed.org/econ/mccracken/fred-databases/>.

where

$$\hat{R}_{i,t+1|t} = \hat{\beta}_{0i} + \hat{\beta}'_{1i} \hat{\mathbf{R}}_t + \dots + \hat{\beta}'_{pi} \hat{\mathbf{R}}_{t-p+1}, \quad i = 1, \dots, n,$$

$\hat{\beta}_{0i}, \hat{\beta}_{1i}, \dots, \hat{\beta}_{pi}$, $i = 1, \dots, n$, are LASSO estimates, $\hat{\mathbf{R}}_t = (\hat{R}_{1,t}, \dots, \hat{R}_{n,t})$, and finally $\hat{R}_{i,t} = Y_{i,t} - \hat{Y}_{i,t|t-1}^{(\text{AR})}$, $i = 1, \dots, n$. The parameters are estimated equation-wise for each one of the 119 variables in the dataset. The penalty parameter is selected by BIC as discussed in Section 4.

3. AR + Principal Component Regression (PCR):

$$\hat{Y}_{i,t+1|t}^{(\text{PCR})} = \hat{Y}_{i,t+1|t}^{(\text{AR})} + \hat{\lambda}'_i \hat{\mathbf{F}}_t,$$

where $\hat{\mathbf{F}}_t$ is the estimate of the $(k \times 1)$ vector of factors \mathbf{F}_t given by principal component analysis of $\hat{\mathbf{R}}_t$, the residuals of the first-stage regression. The parameter λ_i is computed by OLS regression of $\hat{R}_{i,t}$ on $\hat{\mathbf{F}}_t$ in the in-sample window.

4. AR + Full Information (FarmPredict):

$$\hat{Y}_{i,t+1|t}^{(\text{FarmPredict})} = \hat{Y}_{i,t+1|t}^{(\text{PCR})} + \hat{U}_{i,t+1|t}$$

where

$$\hat{U}_{i,t+1|t} = \hat{\theta}_{0i} + \hat{\theta}'_{1i} \hat{\mathbf{U}}_t + \dots + \hat{\theta}'_{pi} \hat{\mathbf{U}}_{t-p+1},$$

$\hat{\mathbf{U}}_t = (\hat{U}_{1,t}, \dots, \hat{U}_{n,t})'$ and $\hat{U}_{i,t} = Y_{i,t} - \hat{Y}_{i,t|t-1}^{(\text{PCR})}$, $i = 1, \dots, n$. The estimates $\hat{\theta}_{0i}, \hat{\theta}_{1i}, \dots, \hat{\theta}_{pi}$, $i = 1, \dots, n$, are given by LASSO.

The forecasts are based on a rolling-window framework of fixed length of 480 observations, starting in January 1960. Therefore, the forecasts start on January 1990. The last forecasts are for December 2019. Note that the AR model only considers information concerning the own past of the variable of interest. SR and PCR expand the information by two opposing routes. While SR uses a sparse combination of the set of variables, PCR considers only a factor structure (dense model). FarmPredict combines these two approaches and uses the full information available.

6.2.3 Brief In-Sample Analysis

We start by looking at the full sample in order to analyse the structure of dependence among the many series considered. We first estimate an autoregressive model of order 4, $\text{AR}(p)$, for each

transformed series. Figure 7 reports the empirical distribution of the OLS estimators of the AR coefficients. Figure 8 shows the distribution of the absolute value of the sum of the estimates. This gives an idea of the persistence of each series. Although we report here the results for AR models of pre-specified order equal to four, in the Supplementary Material we present results for optimal lag selection via the BIC. Only one series has estimated persistence above one. This is the case for *NONBORRES: Reserves of Depository Institutions*, which belongs to group (v): Money and Credit. The reason for such high persistence is due to a major structural break present in the second half of the series. However, 82.35% of the series have estimated persistence below 0.9.¹⁰

We continue by estimating the number of factors when the full sample is used for PCA. We consider two different situations. In the first, we do not include any lag in the basket of variables used to compute the factors. In the second approach, we include four lags of each variable as well. The eigenvalue ratio procedure selects either two (no lags) or a single factor (with lags). The four information criteria of Bai and Ng (2002) as described in Section 5, estimate respectively for the case with no lags (with lags) the following number of factors: six (one), five (one), nine (one), and one (one). Note that the factors are estimated for the residuals of the first-step AR filter. If we remove the *NONBORRES* variable from the sample the results do not change for the eigenvalue ratio procedure. On the other hand, the new numbers of factors selected by the information criteria are as follows: seven (one), six (one), eleven (one), and one (one).

Finally, we apply the testing approach developed in this paper to check for remaining (partial) covariance structure in the data. The tests strongly reject the null of a diagonal matrix when applied to the residuals either of the first or the second stages of the methodology. This serves as evidence that *FarmPredict* may be a useful modeling approach for this macroeconomic dataset.

6.2.4 Forecasting Results

For each of the four models described above, we report a number of performance metrics in Table 6. The table presents the frequency each model has the best performance among the four alternatives. Numbers between parentheses indicates the frequency each model is the second, third, and fourth best. We report the results for each one of the eight sectors as well as for the set of all 119 variables. We show the results for two methods to determine the number of factors. Panel (a) reports the results for the eigenvalue ratio method while Panel (d) presents the results for the information criterion IC_4 . Criteria IC_1 , IC_2 , and IC_3 select a very large number of factors and we relegate them

¹⁰Conventional unit-root tests also reject the null of unit-root for all but one of the series.

to the supplementary material. Panels (c) and (d) in the table show the results for the cases where the number of factors are kept fixed ($r = 1$ or $r = 2$).

FarmPredict is the model which is ranked first more frequently when all the series are considered. It is also the best model for the following groups: output and income, labor market, housing, and consumption, orders and inventories. The **AR** model is best for the following groups: money and credit and stock market. The sparse regression is superior also for two groups: interest and exchange rates and prices.

7 Conclusions

In this paper we propose a new methodology which bridges the gap between sparse regressions and factor models and evaluate the gains of increasing the information set via factor augmentation. Our proposal consists in several steps. In the first one, we filter the data for known factors (trends, seasonal adjustments, covariates). In the second step, we estimate a latent factor structure. Finally, in the last part of the procedure we estimate a sparse regression for the idiosyncratic components. Furthermore, we also propose a new test for remaining structure in both high-dimensional covariance and partial covariance matrices. Our test can be used to evaluate the benefits of adding more structure in the model. Our paper has also a number of important side results. First, we proved consistency of kernel estimation of long-run covariance matrices in high-dimensions where both the number of observations and variables grows. Second, we derive the theoretical properties of factor estimation on the residuals of a first step process. Third, the proposed test can be used as a diagnostic tool for factor models.

We evaluate our methodology with both simulations and real data. The simulations show the test has good size and power properties even when the true number of factors is unknown and must be determined from the data. However, if the number of factors is underestimated, we observe size distortions. This is specially the case when the eigenvalue ratio test is used to determine the number of latent factors. The simulations also show that there are major informational gains when combining factor models and sparse regressions in a forecasting exercise. Two applications are considered in the paper.

A Proof of the Theorems

Throughout the proofs we use the equivalence

$$\|X\|_{\psi_p} < \infty \iff \mathbb{P}(|X| > x) = O(\psi_p^{-1}(x)) \quad \text{as } x \rightarrow \infty,$$

for any random variable X and $\psi_p \in \Psi$, combined with Lemma 6 in Carvalho et al. (2018) and Lemma 1 in Masini and Medeiros (2019). The key ingredients of the lemmas are a Marcinkiewicz-Zygmund type inequality for strong mixing sequences to deal with the polynomial tails (Rio, 1994; Doukhan and Louhichi, 1999) and a Bernstein inequality under strong mixing conditions to control exponential tails (Merlevède et al. (2009) - Theorem 2).

A.1 Proof of Theorem 1

We first upper bound $\|\hat{R}_{it} - R_{it}\|_{\psi}$. By subsequent application of Hölder's inequality we have

$$\begin{aligned} |\hat{R}_{it} - R_{it}| &= |(\hat{\gamma}_i - \gamma_i)' \mathbf{W}_{it}| \\ &\leq \|\hat{\gamma}_i - \gamma_i\|_1 \|\mathbf{W}_{it}\|_{\infty} \\ &= \|\hat{\Sigma}_i^{-1} \hat{\mathbf{v}}_i\|_1 \|\mathbf{W}_{it}\|_{\infty} \\ &\leq k^2 \|\hat{\Sigma}_i^{-1}\|_{\max} \|\hat{\mathbf{v}}_i\|_{\infty} \|\mathbf{W}_{it}\|_{\infty}, \end{aligned}$$

where $\hat{\Sigma}_i := \mathbf{W}'_i \mathbf{W}_i / T$ and $\hat{\mathbf{v}}_i := \mathbf{W}'_i \mathbf{U}_i / T$. Then by the Cauchy-Schwartz conjugate

$$\|\hat{R}_{it} - R_{it}\|_{\psi_{p/4}} \leq k^2 \|\|\hat{\Sigma}_i^{-1}\|_{\max}\|_{\psi_p} \|\|\hat{\mathbf{v}}_i\|_{\infty}\|_{\psi_{p/2}} \|\|\mathbf{W}_{it}\|_{\infty}\|_{\psi_p}.$$

The first term is bounded by Assumption 3(b). For the second term we have: $\|W_{it\ell} U_{it}\|_{\psi_{p/2}} \leq \|W_{it\ell}\|_{\psi_p} \|U_{it}\|_{\psi_p} \leq C^2$ by Assumption 3(a). Then, $\{W_{it\ell} U_{it}\}_{t>0}$ is a zero-mean strong mixing with exponential decay sequence (Assumption 3(c)) with bounded $\psi_{p/2}$ -norm. Therefore, $\|\|\hat{\mathbf{v}}_i\|_{\infty}\|_{\psi_{p/2}} = O(1/\sqrt{T})$ uniformly in $i \leq n$. Finally, the last term is bounded by the maximal inequality (van der Vaart and Wellner (1996) - Lemma 2.2.2) and Assumption 3(a). The first result follows.

A.2 Proof of Theorem 2

The proof is an adaption of the proof of Theorem 4 and Corollary 1 in Fan et al. (2013), henceforth FLM, to include the estimation error in the sample covariance matrix. For part (a), we pick up from expression (A.1) in Bai (2003) to obtain the following identity

$$\hat{\mathbf{f}}_t - \mathbf{H}\mathbf{F}_t = \left(\frac{\mathbf{V}}{n}\right)^{-1} \left[\frac{1}{T} \sum_{s=1}^T \hat{\mathbf{f}}_s \frac{\mathbb{E}(\mathbf{U}'_s \mathbf{U}_t)}{n} + \frac{1}{T} \sum_{s=1}^T \hat{\mathbf{f}}_s \tilde{\zeta}_{st} + \frac{1}{T} \sum_{s=1}^T \hat{\mathbf{f}}_s \tilde{\eta}_{st} + \frac{1}{T} \sum_{s=1}^T \hat{\mathbf{f}}_s \tilde{\xi}_{st} \right], \quad (\text{A.1})$$

where $\tilde{\zeta}_{st}$, $\tilde{\eta}_{st}$ and $\tilde{\xi}_{st}$ are defined before Lemma B.3.

By Assumptions 2(d) and 3(a) and the maximal inequality we have $\|\mathbf{R}\|_{\max} \leq r \|\mathbf{A}\|_{\max} \|\mathbf{F}\|_{\max} + \|\mathbf{U}\|_{\max} = O_P(\psi^{-1}(nT))$. Applying Lemma B.14 we conclude that $\|\hat{\Sigma} - \tilde{\Sigma}\|_{\max} = O_P(\omega(\psi^{-1}(nT) + \omega)) = O_P(1)$, where the last assumption by the Theorem assumption. Finally $\psi^{-1}(n^2)/\sqrt{T} = O(1)$ also by assumption then $\|\frac{\mathbf{V}}{n}\|^{-1} = O_P(1)$ by Lemma B.6. Using the results (a)-(d) of Lemma B.5 we can bound in probability each of the terms in brackets of (A.1) in ℓ_2 norm uniformly in $t \leq T$ and obtain the result (a).

For part (b) we use the fact that $\hat{\Lambda} := \hat{\mathbf{R}}\hat{\mathbf{F}}/T$ and the normalization $\hat{\mathbf{F}}'\hat{\mathbf{F}} = \mathbf{I}_r$ to write

$$\hat{\lambda}_i - \mathbf{H}\lambda_i = \frac{1}{T} \sum_{t=1}^T \mathbf{H}\mathbf{F}_t \tilde{U}_{it} + \frac{1}{T} \sum_{t=1}^T \hat{R}_{it} (\hat{\mathbf{F}}_t - \mathbf{H}\mathbf{F}_t) + \mathbf{H} \left(\frac{1}{T} \sum_{t=1}^T \mathbf{F}_t \mathbf{F}'_t - \mathbf{I}_r \right) \lambda_i. \quad (\text{A.2})$$

The first term can be upper bounded in ℓ_2 norm uniformly in $i \leq n$ by

$$\sqrt{r} \|\mathbf{H}\| \max_{i \leq n} \max_{j \leq r} \left| \frac{1}{T} \sum_{t=1}^T F_{jt} \tilde{U}_{it} \right| = O_P(1) O_P[\psi_{p/2}^{-1}(n)/\sqrt{T} + \omega],$$

where the equality follows from Lemma B.6(b) and (e). The ℓ_2 norm of the second term is upper bounded uniformly in $i \leq n$ by

$$\left(\max_{i \leq n} \frac{1}{T} \sum_{t=1}^T \hat{R}_{it}^2 \frac{1}{T} \sum_{t=1}^T \|\hat{\mathbf{F}}_t - \mathbf{H}\mathbf{F}_t\|^2 \right)^{1/2} = \left[O_P(1) O_P\left(\frac{1}{T} + (1/\sqrt{n} + \omega)^2\right) \right]^{1/2},$$

where the first term after the equality follows from Lemma B.6(d) together with the Theorem assumption and the second term from Lemma B.4(e). Finally the last term of (A.2) is upper bounded by

$$\|\mathbf{H}\| \max_{i \leq n} \|\lambda_i\| \left\| \frac{1}{T} \sum_{t=1}^T \mathbf{F}_t \mathbf{F}'_t - \mathbf{I}_r \right\| = O_P(1) O(1) O_P(1/\sqrt{T}),$$

where the last term is $O_P(1/\sqrt{T})$ by the maximum inequality and Assumption 3 Plug the last three displays back into (A.2) yields result (b).

For part (c) we use we have $\|\widehat{\mathbf{U}} - \mathbf{U}\|_{\max} = \|\boldsymbol{\Lambda}\mathbf{F}' - \widehat{\boldsymbol{\Lambda}}\widehat{\mathbf{F}}' + \widehat{\mathbf{R}} - \mathbf{R}\|_{\max} \leq \|\widehat{\boldsymbol{\Lambda}}\widehat{\mathbf{F}}' - \boldsymbol{\Lambda}\mathbf{F}'\|_{\max} + \|\widehat{\mathbf{R}} - \mathbf{R}\|_{\max}$. The last term is $O_P(\omega)$ by assumption. For the first term we use the decomposition

$$\begin{aligned} \widehat{\boldsymbol{\lambda}}_i' \widehat{\mathbf{F}}_t - \boldsymbol{\lambda}_i' \mathbf{F}_t &= (\widehat{\boldsymbol{\lambda}}_i - \mathbf{H}\boldsymbol{\lambda}_i)'(\widehat{\mathbf{F}}_t - \mathbf{H}\mathbf{F}_t) + (\mathbf{H}\boldsymbol{\lambda}_i)'(\widehat{\mathbf{F}}_t - \mathbf{H}\mathbf{F}_t) \\ &\quad + (\widehat{\boldsymbol{\lambda}}_i - \mathbf{H}\boldsymbol{\lambda}_i)' \mathbf{H}\mathbf{F}_t + \boldsymbol{\lambda}_i' (\mathbf{H}'\mathbf{H} - \mathbf{I}_r) \mathbf{F}_t. \end{aligned} \quad (\text{A.3})$$

Therefore, we can upper bound the left hand side as

$$\begin{aligned} |\widehat{\boldsymbol{\lambda}}_i' \widehat{\mathbf{F}}_t - \boldsymbol{\lambda}_i' \mathbf{F}_t| &\leq \|\widehat{\boldsymbol{\lambda}}_i - \mathbf{H}\boldsymbol{\lambda}_i\| \|\widehat{\mathbf{F}}_t - \mathbf{H}\mathbf{F}_t\| + \|\mathbf{H}\boldsymbol{\lambda}_i\| \|\widehat{\mathbf{F}}_t - \mathbf{H}\mathbf{F}_t\| \\ &\quad + \|\widehat{\boldsymbol{\lambda}}_i - \mathbf{H}\boldsymbol{\lambda}_i\| \|\mathbf{H}\mathbf{F}_t\| + \|\boldsymbol{\lambda}_i\| \|\mathbf{F}_t\| \|\mathbf{H}'\mathbf{H} - \mathbf{I}_r\|. \end{aligned}$$

Now we bound in probability each of the four term above uniformly in $i \leq n$ and $t \leq T$. The first one is given by part (a) and (b). $\max_{i \leq n} \|\mathbf{H}\boldsymbol{\lambda}_i\| \leq \|\mathbf{H}\| \max_{i \leq n} \|\boldsymbol{\lambda}_i\| \leq O_P(1)r\|\boldsymbol{\Lambda}\|_{\max} = O_P(1)$ by Lemma B.6(b) and Assumption 2(d), thus the second term is bounded by part (a). Similarly for the third term $\max_{t \leq T} \|\mathbf{H}\mathbf{F}_t\| \leq \|\mathbf{H}\| \max_{t \leq T} \|\mathbf{F}_t\| = O_P(1)O_P(\psi^{-1}(T)) = O_P(\psi^{-1}(T))$ by Lemma B.6(b) and Assumption 2(a). Finally $\|\mathbf{H}'\mathbf{H} - \mathbf{I}_r\| = O_P(1/\sqrt{T} + 1/\sqrt{n} + \omega)$ by Lemma B.6(c) hence the last term is $O_P[\psi^{-1}(T)(1/\sqrt{T} + 1/\sqrt{n} + \omega)]$ by Assumptions 2(d) and 3(a).

A.3 Proof of Theorem 3

We have that $L(\widehat{\boldsymbol{\theta}}_\xi) + \xi\|\widehat{\boldsymbol{\theta}}_\xi\|_1 \leq L(\boldsymbol{\theta}) + \xi\|\boldsymbol{\theta}\|_1$ for all $\boldsymbol{\theta} \in \mathbb{R}^n$ by definition of $\widehat{\boldsymbol{\theta}}_\xi$, where $L(\boldsymbol{\theta}) := \|\widehat{\mathbf{u}}_y - \boldsymbol{\theta}'\widehat{\mathbf{U}}_x\|_2^2/T$. Also, since $L(\boldsymbol{\theta})$ is a quadratic function, it implies that $(\widehat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta})'\nabla^2 L(\boldsymbol{\theta})(\widehat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta}) \leq -\nabla L(\boldsymbol{\theta})'(\widehat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta}) + \xi(\|\boldsymbol{\theta}\|_1 - \|\widehat{\boldsymbol{\theta}}_\xi\|_1)$. By Holder's inequality we have $|\nabla L(\boldsymbol{\theta})'(\widehat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta})| \leq \|\nabla L(\boldsymbol{\theta})\|_\infty \|\widehat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta}\|_1$ and by assumption $\xi \geq 2\|\nabla L(\boldsymbol{\theta})\|_\infty$ then we have

$$(\widehat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta})'\nabla^2 L(\boldsymbol{\theta})(\widehat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta}) \leq \xi/2\|\widehat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta}\|_1 + \xi(\|\boldsymbol{\theta}\|_1 - \|\widehat{\boldsymbol{\theta}}_\xi\|_1). \quad (\text{A.4})$$

For any index set $\mathcal{S} \in [n]$, by the decomposability of the ℓ_1 norm (refer to Definition 1 in Negahban et al. (2012)) followed by the triangle inequality we have $\|\widehat{\boldsymbol{\theta}}_\xi\|_1 = \|\widehat{\boldsymbol{\theta}}_{\xi, \mathcal{S}}\|_1 + \|\widehat{\boldsymbol{\theta}}_{\xi, \mathcal{S}^c}\|_1 \geq \|\boldsymbol{\theta}_{\mathcal{S}}\|_1 - \|\widehat{\boldsymbol{\theta}}_{\xi, \mathcal{S}} - \boldsymbol{\theta}_{\mathcal{S}}\|_1 + \|\widehat{\boldsymbol{\theta}}_{\xi, \mathcal{S}^c}\|_1$ and $\|\widehat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta}\|_1 = \|\widehat{\boldsymbol{\theta}}_{\xi, \mathcal{S}} - \boldsymbol{\theta}_{\mathcal{S}}\|_1 + \|\widehat{\boldsymbol{\theta}}_{\xi, \mathcal{S}^c} - \boldsymbol{\theta}_{\mathcal{S}^c}\|_1 \leq \|\widehat{\boldsymbol{\theta}}_{\xi, \mathcal{S}} - \boldsymbol{\theta}_{\mathcal{S}}\|_1 + \|\widehat{\boldsymbol{\theta}}_{\xi, \mathcal{S}^c} - \boldsymbol{\theta}_{\mathcal{S}^c}\|_1$. Plugging

it back in (A.4) yields

$$2(\hat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta})' \nabla^2 L(\boldsymbol{\theta})(\hat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta}) + \xi \|\hat{\boldsymbol{\theta}}_{\xi, \mathcal{S}^c} - \boldsymbol{\theta}_{\mathcal{S}^c}\|_1 \leq 3\xi \|\hat{\boldsymbol{\theta}}_{\xi, \mathcal{S}} - \boldsymbol{\theta}_{\mathcal{S}}\|_1 + 4\xi \|\boldsymbol{\theta}_{\mathcal{S}^c}\|_1. \quad (\text{A.5})$$

We then conclude that any minimizer $\hat{\boldsymbol{\theta}}_\xi$ of (3.6) and $\boldsymbol{\theta} \in \mathbb{R}^n$ obeys $\hat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta} \in \mathbb{C}(\mathcal{S}, \boldsymbol{\theta}) := \{\boldsymbol{x} \in \mathbb{R}^n : \|\boldsymbol{x}_{\mathcal{S}^c}\|_1 \leq 3\|\boldsymbol{x}_{\mathcal{S}}\|_1 + 4\|\boldsymbol{\theta}_{\mathcal{S}^c}\|_1\}$. If we take $\boldsymbol{\theta} = \boldsymbol{\theta}_0$ and $\mathcal{S} = \mathcal{S}_0 := \{i : \theta_{0,i} \neq 0\}$ then $\hat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta}_0 \in \mathbb{C}_0 := \mathbb{C}(\mathcal{S}_0, \boldsymbol{\theta}_0)$. Note that \mathbb{C}_0 is a cone in \mathbb{R}^n that does not depend on $\boldsymbol{\theta}_0$ as $\|\boldsymbol{\theta}_{0, \mathcal{S}_0^c}\| = 0$. Moreover by definition of the compatibility constant $\kappa := \kappa(\hat{\mathbf{U}}_x \hat{\mathbf{U}}_x / T, \mathcal{S}_0, 3)$ we have that $\|\hat{\boldsymbol{\theta}}_{\xi, \mathcal{S}} - \boldsymbol{\theta}_{\mathcal{S}}\|_1 \leq (\hat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta})' \nabla^2 L(\boldsymbol{\theta})(\hat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta}) \sqrt{|\mathcal{S}_0|} / \kappa$. Apply this inequality (A.5) and use the fact, $4ab < a^2 + 4b^2$ for non-negative $a, b \in \mathbb{R}$ to obtain

$$(\hat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta})' \nabla^2 L(\boldsymbol{\theta})(\hat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta}) + \xi \|\hat{\boldsymbol{\theta}}_\xi - \boldsymbol{\theta}\|_1 \leq 4\xi^2 |\mathcal{S}_0| / \kappa. \quad (\text{A.6})$$

Finally, we have by assumption $\|\hat{\mathbf{U}} - \mathbf{U}\|_{\max} \leq C_1$, $\|\mathbf{U}\|_{\max} \leq C_2$ and $C_1(2C_2 + C_1) \leq \frac{\kappa_0}{32|\mathcal{S}|}$ which, in turn fulfills the assumptions of Lemma B.14 with $\zeta = 3$ and $\alpha = 1/2$. Therefore, we conclude that $\kappa \geq \kappa_0/2$ and we have the result.

A.4 Proof of Theorem 4

We use in this proof the following additional notation for short: For every random vector \mathbf{X} , we denote by $\boldsymbol{\Sigma}_{\mathbf{X}}$ its covariance matrix, $d_{\mathbf{X}}$ the diagonal of $\boldsymbol{\Sigma}_{\mathbf{X}}$ and $\sigma_{\mathbf{X}}^2 := \|d_{\mathbf{X}}\|_{\infty}$. Also, \mathbf{X}_G denotes zero-mean Gaussian random vector defined in the same probability space, independent of \mathbf{X} and with the same covariance matrix of \mathbf{X} . Finally, for every pair of random vectors \mathbf{X}, \mathbf{Y} of the same dimension and scalar $s > 0$ define

$$\rho(\mathbf{X}, \mathbf{Y}) := \sup_{t \in \mathbb{R}} |\mathbb{P}(\|\mathbf{X}\|_{\infty} \leq t) - \mathbb{P}(\|\mathbf{Y}\|_{\infty} \leq t)|$$

$$\Delta(\mathbf{X}, s) := \sup_{t \in \mathbb{R}} \mathbb{P}(t \leq \|\mathbf{X}\|_{\infty} \leq t + s)$$

Combining equations (83)–(86) in Giessing and Fan (2020) gives us the following basic inequality

$$\begin{aligned} |\mathbb{P}(S \leq c^*(\tau)) - \tau| &\leq \rho(\tilde{\mathbf{Q}}, \tilde{\mathbf{Q}}_G) + \inf_{\delta_1 > 0} \left\{ \sqrt{\delta_1} \frac{\log n}{\omega_{\max}} + \mathbb{P}(\|\hat{\mathbf{Y}} - \mathbf{Y}\|_{\max} > \delta_1) \right\} \\ &\quad + \inf_{\delta_2 > 0} \left\{ \delta_2 \frac{\sqrt{\log n}}{\omega_{\max}} + \mathbb{P}(\|\mathbf{Q} - \tilde{\mathbf{Q}}\|_{\infty} > \delta_2) \right\} \end{aligned} \quad (\text{A.7})$$

where $\tilde{\mathbf{Q}}$ is defined below.

We start by Bounding the first term to the right-hand side of (A.7). Here we adapt the classical "big block-small block" technique proposed by Bernstein in the context of proving CLT under mixing conditions, which was also used in the proof of Theorem E.1 in Chernozhukov et al. (2018). Consider two sequences of non-negative integers $a := a_T$ and $b := b_T$ such that $b < a$, $a + b \leq T$, $\min\{a, b\} \rightarrow \infty$, $a = o(T)$ and $b = o(a)$ as $T \rightarrow \infty$. Let $m := \lfloor T/(a+b) \rfloor$ and define for $j \in \{1, \dots, m\}$ consecutive blocks of size a and b with index set $\mathcal{A}_j := \{(j-1)(a+b) + 1, \dots, (j-1)(a+b) + a\}$ and $\mathcal{B}_j := \{(j-1)(a+b) + a + 1, \dots, j(a+b)\}$. Finally set $\mathcal{C} := \{m(a+b) + 1, \dots, T\}$, which might be empty.

$$\mathbf{A}_j := \frac{1}{\sqrt{a}} \sum_{t \in \mathcal{A}_j} \tilde{\mathbf{D}}_t \quad \mathbf{B}_j = \frac{1}{\sqrt{b}} \sum_{t \in \mathcal{B}_j} \tilde{\mathbf{D}}_t; \quad \mathbf{C} = \frac{1}{\sqrt{|\mathcal{C}|}} \sum_{t \in \mathcal{C}} \tilde{\mathbf{D}}_t,$$

such that

$$\tilde{\mathbf{Q}} := \frac{1}{\sqrt{T}} \sum_{t=1}^T \tilde{\mathbf{D}}_t = \underbrace{\sqrt{\frac{ma}{T}} \left(\frac{1}{\sqrt{m}} \sum_{j=1}^m \mathbf{A}_j \right)}_{=: \mathbf{V}} + \underbrace{\sqrt{\frac{mb}{T}} \left(\frac{1}{\sqrt{m}} \sum_{j=1}^m \mathbf{B}_j \right)}_{=: \mathbf{L}} + \sqrt{\frac{T - m(a+b)}{T}} \mathbf{C}$$

Now let $\tilde{\mathbf{V}} := \frac{1}{\sqrt{m}} \sum_{j=1}^m \tilde{\mathbf{A}}_j$ where $\{\tilde{\mathbf{A}}_t, 1 \leq t \leq m\}$ is an independent sequence such that \mathbf{A}_t and $\tilde{\mathbf{A}}_t$ have the same distribution for all $1 \leq t \leq m$. Similarly define $\tilde{\mathbf{L}} := \frac{1}{\sqrt{m}} \sum_{j=1}^m \tilde{\mathbf{B}}_j$. Lemma B.7 give us for any scalar $s > 0$

$$\begin{aligned} \rho(\tilde{\mathbf{Q}}, \tilde{\mathbf{Q}}_G) &\leq \rho(\tilde{\mathbf{V}}, \tilde{\mathbf{V}}_G) + \rho\left(\sqrt{\frac{ma}{T}} \tilde{\mathbf{V}}_G, \tilde{\mathbf{Q}}_G\right) + \Delta\left(\sqrt{\frac{ma}{T}} \tilde{\mathbf{V}}_G, s\right) \\ &\quad + \mathbb{P}\left(\sqrt{\frac{mb}{T}} \|\tilde{\mathbf{L}}\|_\infty > s\right) + \rho(\mathbf{V}, \tilde{\mathbf{V}}) + \rho(\mathbf{L}, \tilde{\mathbf{L}}). \end{aligned} \quad (\text{A.8})$$

Notice that for any measurable $A \subseteq \mathbb{R}^2$ we have $|\mathbb{P}[(\mathbf{A}_1, \mathbf{A}_2) \in A] - \mathbb{P}[\tilde{\mathbf{A}}_1, \tilde{\mathbf{A}}_2,]| \leq \alpha_b$ where $\{\alpha_n, n \in \mathbb{N}\}$ denote the α -mixing coefficient of the sequence $(\tilde{\mathbf{D}}_t)$ which is the same of the sequence (\mathbf{U}_t) . Then the last two terms in (A.8) can be upper bounded by $(m-1)\alpha_b$ and $(m-1)\alpha_a$ respectively by induction. Since α_n is non-increasing in n and $a \geq b$ we have that

$$\rho(\mathbf{V}, \tilde{\mathbf{V}}) + \rho(\mathbf{L}, \tilde{\mathbf{L}}) \leq 2(m-1)\alpha_b \leq 2T \exp(-2cb). \quad (\text{A.9})$$

where we use Assumption 3(c) to obtain the last inequality.

For the fourth term we have by the maximal inequality followed by Markov's inequality $\mathbb{P}(\sqrt{\frac{mb}{T}} \|\tilde{\mathbf{L}}\|_\infty >$

$s) \leq \left[\psi \left(\frac{s\sqrt{T}}{C_\psi \psi_{p/2}^{-1}(n) \sqrt{mb}} \right) \right]^{-1}$ and the anti-concentration inequality for Gaussian random vectors (Theorem 7 in Giessing and Fan (2020) with $p = \infty$) $\Delta(\sqrt{\frac{ma}{T}} \tilde{\mathbf{V}}_G, s) \lesssim \frac{Ts\sqrt{\log n}}{ma\sigma_{\tilde{\mathbf{V}}}}$. Set $s = \frac{C_\psi \psi_{p/2}^{-1}(n) \sqrt{mb}}{\sqrt{T}} \psi_{p/2}^{-1}(T^\gamma)$ for some $\gamma > 0$ then

$$\Delta\left(\sqrt{\frac{ma}{T}} \tilde{\mathbf{V}}_G, s\right) + \mathbb{P}\left(\sqrt{\frac{mb}{T}} \|\tilde{\mathbf{L}}\|_\infty > s\right) \lesssim \frac{T}{ma} \sqrt{\frac{mb}{T}} \frac{\sqrt{\log n} \psi_{p/2}^{-1}(n) \psi_{p/2}^{-1}(T^\gamma)}{\sigma_{\tilde{\mathbf{V}}}} + \frac{1}{T^\gamma} \quad (\text{A.10})$$

For the second term we have from Rio (2013) that, for every $\epsilon > 0$,

$$|[\mathbf{M}_\ell]_{ij}| = |\text{Cov}(\tilde{D}_{it}, \tilde{D}_{j,t-\ell})| \leq 2\alpha_\ell^{\epsilon/(2+\epsilon)} \|\tilde{D}_{it}\|_{2+\epsilon} \|\tilde{D}_{j,t-\ell}\|_{2+\epsilon}.$$

Hence, from Assumption 3 we have that $\|\mathbf{M}_\ell\|_{\max} \lesssim \exp(-2c\frac{\epsilon}{2+\epsilon}\ell)$ and

$$\begin{aligned} \left\| \frac{ma}{T} \Sigma_{\tilde{\mathbf{V}}_G} - \Sigma_{\tilde{\mathbf{Q}}_G} \right\|_{\max} &\leq \left(1 - \frac{ma}{T}\right) \|\Sigma_{\tilde{\mathbf{V}}}\|_{\max} + \|\Sigma_{\tilde{\mathbf{Q}}} - \Sigma_{\tilde{\mathbf{V}}}\|_{\max} \\ &\leq \left(\frac{b}{a+b} + \frac{a}{T}\right) \|\Sigma_{\tilde{\mathbf{V}}}\|_{\max} + \frac{1}{a} \sum_{|\ell|<a} |\ell| \|\mathbf{M}_\ell\|_{\max} + \sum_{a \leq |\ell| < T} \|\mathbf{M}_\ell\|_{\max} \\ &\lesssim \frac{b}{a} + \frac{a}{T} + \frac{1}{a} + T \exp(-2c\frac{\epsilon}{2+\epsilon}a), \end{aligned}$$

where we use the fact that $\Sigma_{\tilde{\mathbf{V}}_G} = \Sigma_{\tilde{\mathbf{V}}} = \Sigma_{\tilde{\mathbf{A}}_j} = \Sigma_{\mathbf{A}_j} = \sum_{|\ell|<a} (1 - |\ell|/a) \mathbf{M}_\ell$, $\Sigma_{\tilde{\mathbf{Q}}_G} = \Sigma_{\tilde{\mathbf{Q}}} = \sum_{|\ell|<T} (1 - |\ell|/T) \mathbf{M}_\ell$, $\sum_{|\ell|<a} |\ell| \|\mathbf{M}_\ell\|_{\max} \leq c$ for some $c < \infty$ and $\sum_{a \leq |\ell| < T} \|\mathbf{M}_\ell\|_{\max} \lesssim T \exp(-2c\frac{\epsilon}{2+\epsilon}a)$. Finally, we can bound the second term using Theorem 8 in Giessing and Fan (2020). In particular for $p = \infty$ it implies that

$$\rho\left(\sqrt{\frac{ma}{T}} \tilde{\mathbf{V}}_G, \tilde{\mathbf{Q}}_G\right) \lesssim \frac{\log n \sqrt{\left\| \frac{ma}{T} \Sigma_{\tilde{\mathbf{V}}_G} - \Sigma_{\tilde{\mathbf{Q}}_G} \right\|_{\max}}}{\sqrt{\frac{ma}{T}} \sigma_{\tilde{\mathbf{V}}} \vee \sigma_{\tilde{\mathbf{Q}}}} \lesssim \frac{\sqrt{\frac{T}{ma}} \log n \sqrt{\frac{b}{a} + \frac{a}{T} + \frac{1}{a} + T \exp(-\frac{2c\epsilon}{2+\epsilon}a)}}{\sigma_{\tilde{\mathbf{V}}} \vee \sigma_{\tilde{\mathbf{Q}}}} \quad (\text{A.11})$$

For the first term we have that $\|\tilde{D}_{it}\|_{\psi_{p/2}}$ is uniformly (upper) bounded by Assumption 3(a) then so is $\|\tilde{A}_{it}\|_{\psi_{p/2}} = \|A_{it}\|_{\psi_{p/2}} = \|\frac{1}{\sqrt{a}} \sum_{s \in \mathcal{A}_t} \tilde{D}_{is}\|_{\psi_{p/2}}$. Also $(\mathbb{E}(\max_i |\tilde{A}_{it}|)^3)^{1/3} \lesssim \|\max_i |\tilde{A}_{it}|\|_{\psi_{p/2}} \lesssim \psi_{p/2}^{-1}(n) \max_i \|\tilde{A}_{it}\|_{\psi_{p/2}} \lesssim \psi_{p/2}^{-1}(n)$. Since $\{\tilde{A}_t, 1 \leq t \leq m\}$ is an iid sequence of random vector Theorem 5 in Giessing and Fan (2020) gives us

$$\rho(\tilde{\mathbf{V}}, \tilde{\mathbf{V}}_G) \lesssim \frac{(\log n)^{7/6} \psi_{p/2}^{-1}(n)}{T^{1/6} \sigma_{\tilde{\mathbf{V}}}}. \quad (\text{A.12})$$

By the triangle inequality we have that $\sigma_{\tilde{\mathbf{V}}}^2 \geq \sigma_{\tilde{\mathbf{Q}}}^2 - \|d_{\tilde{\mathbf{Q}}} - d_{\tilde{\mathbf{V}}}\|_{\max} \geq \underline{c} - \|\Sigma_{\tilde{\mathbf{Q}}} - \Sigma_{\tilde{\mathbf{V}}}\|_{\max} \gtrsim$

$c - \frac{1}{a} - T \exp(-2c \frac{\epsilon}{2+\epsilon} a)$. By setting $a = \lceil \sqrt{T} \rceil$ we conclude that $\sigma_{\mathbf{V}}^2$ is eventually bounded away from zero for large enough T . If we further set $b = \lceil \log T/c \rceil$ and $\gamma = 1/4$ and apply (A.9)-(A.12) to bound the right-hand side of (A.7) we obtain

$$\rho(\tilde{\mathbf{Q}}, \tilde{\mathbf{Q}}_G) = O \left[\frac{(\log n)^{7/6} \psi_{p/2}^{-1}(n)}{T^{1/6}} + \frac{\sqrt{\log T \log n} \psi_{p/2}^{-1}(n) \psi_{p/2}^{-1}(T^{1/4})}{T^{1/4}} \right]. \quad (\text{A.13})$$

Finally, we now bound the last two terms appearing in (A.7). Let γ_1 and γ_2 be positive sequences depending on n and T such that $\|\hat{\mathbf{Y}} - \mathbf{Y}\|_{\max} = O_P(\gamma_1)$ and $\|\mathbf{Q} - \tilde{\mathbf{Q}}\|_{\infty} = O_P(\gamma_2)$. Suppose we can state conditions under which

$$\log^3 n (\gamma_1 \vee \gamma_2) = o(1) \quad T, n \rightarrow \infty \quad (\text{A.14})$$

Then we have the the last two terms vanish in probability if we set $\delta_1 = \gamma_1 \log n$ and $\delta_2 = \gamma_2 \log n$ in (A.7). Lemma B.8 and Lemma B.10 give us expressions for γ_1 and γ_2 , respectively, which combined with the rate assumptions in the theorem implies (A.14).

B Additional Lemmas

Lemma B.1. *Let a_j and b_j denote the j -th eigenvalue in decreasing order of $\mathbf{\Sigma}$ and $\mathbf{\Lambda}\mathbf{\Lambda}'$ respectively.*

Then, under Assumption 2(b) and (c):

(a) $b_j \asymp n$ for $1 \leq j \leq r$

(b) $\max_{j \leq n} |a_j - b_j| = O(1)$

(c) $a_j \asymp n$ for $1 \leq j \leq r$.

Proof. Result (a) follows from the fact that the r eigenvalues of $\mathbf{\Lambda}'\mathbf{\Lambda}$ are also (the only r non-zero) eigenvalues of $\mathbf{\Lambda}\mathbf{\Lambda}'$ and Assumption 2(b). Part (b) follows from Weyl's inequality that implies $\max_{j \leq n} |a_j - b_j| \leq \|\mathbf{\Sigma} - \mathbf{\Lambda}\mathbf{\Lambda}'\| = O(1)$, where the last equality follows from Assumption 2(c). Finally result (c) follows from part (a) and (b) and the (reverse) triangle inequality. \square

Recall that $\mathbf{\Sigma}$ be the $(n \times n)$ covariance matrix of $\mathbf{U}_t = \mathbf{Z}_t - \mathbf{\Gamma}\mathbf{W}_t$. Let $\tilde{\mathbf{\Sigma}} := \frac{1}{T} \sum_{t=1}^T \mathbf{U}_t \mathbf{U}_t'$ and $\hat{\mathbf{\Sigma}}$ the same as $\tilde{\mathbf{\Sigma}}$ but with $\mathbf{\Gamma}$ replaced by the estimator $\hat{\mathbf{\Gamma}}$. Also let \hat{a}_j denote the j -th eigenvalue in decreasing order of $\hat{\mathbf{\Sigma}}$

Lemma B.2. Let ω_1 be a non-negative sequence of n and T such that $\|\widehat{\Sigma} - \widetilde{\Sigma}\|_{\max} = O_P(\omega_1)$. Then, under the Assumptions 2 and 3:

- (a) $\|\widehat{\Sigma} - \Sigma\|_{\max} = O_P[\omega_1 + \psi_{p/2}^{-1}(n^2)/\sqrt{T}]$
- (b) $\max_{j \leq n} |\widehat{a}_j - a_j| = O_P[n(\omega_1 + \psi_{p/2}^{-1}(n^2)/\sqrt{T})]$
- (c) $\widehat{a}_j \asymp_P n$ for $j \leq r$ provided that $\omega_1 + \psi_{p/2}^{-1}(n^2)/\sqrt{T} = O_P(1)$

Proof. Part (a) follows by triangle inequality followed by the maximum inequality since $\|\widehat{\Sigma} - \Sigma\|_{\max} \leq \|\widehat{\Sigma} - \widetilde{\Sigma}\|_{\max} + \|\widetilde{\Sigma} - \Sigma\|_{\max} = O_P(\omega_1) + O_P(\psi_{p/2}^{-1}(n^2)/\sqrt{T})$. Part (b) follows from Weyl's inequality, the fact that $\|\widehat{\Sigma} - \Sigma\| \leq n\|\widehat{\Sigma} - \Sigma\|_{\max}$ and part (a). Part (c) follows from the triangle inequality combined with part (b) and Lemma B.1(c). \square

The Lemmas B.3-B.6 below are an adaption of Lemmas 8-10 in Fan et al. (2013), henceforth FLM, to include the estimation error in the sample covariance matrix. To avoid confusion and make it easier for the read to follow through the changes we use the same notation adopted in FLM. In particular, if δ_{it} denotes the (i, t) element of $\Delta := \widehat{\mathbf{R}} - \mathbf{R}$ then $\widetilde{U}_{it} = U_{it} + \delta_{it}$ for $i \in [n]$ and $t \in [T]$. Also, we consider that $\|\Delta\|_{\max} = O_P(\omega)$ for some non-negative sequence ω depending on n and T .

Define:

$$\begin{aligned} \widetilde{\zeta}_{st} &:= \frac{\widetilde{U}'_s \widetilde{U}_t}{n} - \frac{\mathbb{E}(U'_s U_t)}{n} = \left(\frac{U'_s U_t}{n} - \frac{\mathbb{E}(U'_s U_t)}{n} \right) + \left(\frac{U'_s \delta_t}{n} + \frac{\delta'_s U_t}{n} + \frac{\delta'_s \delta_t}{n} \right) =: \zeta_{st} + \zeta_{st}^* \\ \widetilde{\eta}_{st} &:= \frac{\mathbf{f}'_s \sum_{i=1}^n \lambda_i \widetilde{U}_{it}}{n} = \frac{\mathbf{f}'_s \sum_{i=1}^n \lambda_i U_{it}}{n} + \frac{\mathbf{f}'_s \sum_{i=1}^n \lambda_i \delta_{it}}{n} =: \eta_{st} + \eta_{st}^* \\ \widetilde{\xi}_{st} &:= \frac{\mathbf{F}'_t \sum_{i=1}^n \lambda_i \widetilde{U}_{is}}{n} = \frac{\mathbf{F}'_t \sum_{i=1}^n \lambda_i U_{is}}{n} + \frac{\mathbf{F}'_t \sum_{i=1}^n \lambda_i \delta_{is}}{n} = \xi_{st} + \xi_{st}^*. \end{aligned}$$

Lemma B.3. Under Assumption 3:

- (a) $\zeta_{st} = O_P(1/\sqrt{n})$
- (b) $\eta_{st} = O_P(1/\sqrt{n})$
- (c) $\xi_{st} = O_P(1/\sqrt{n})$
- (d) $\zeta_{st}^* = O_P(\omega + \omega^2)$ and $\max_{s,t \leq T} \zeta_{st}^* = O_P(\psi^{-1}(nT)\omega + \omega^2)$
- (e) $\eta_{st}^* = O_P(\omega)$
- (f) $\xi_{st}^* = O_P(\omega)$.

Proof. Parts (a), (b) and (c) are straightforward. For (d) we have that $\frac{1}{n}\mathbf{U}'_s\mathbf{U}_t = O_P(1)$ and $\frac{1}{n}\boldsymbol{\delta}'_s\boldsymbol{\delta}_t \leq \|\boldsymbol{\Delta}\|_{\max}^2 = O_P(\omega^2)$ then the other two terms in parentheses in the definition of ζ_{st}^* are $O_P(\omega)$ by the Cauchy-Schwartz inequality. Part (e) and (f) follows by similar arguments.

$$\max_{t \leq T} \frac{1}{T} \sum_{s=1}^T \left(\frac{1}{n} \boldsymbol{\delta}'_s \mathbf{U}_t \right)^2 = \max_{t \leq T} \frac{1}{n^2} \mathbf{U}'_t \left(\frac{1}{T} \sum_{s=1}^T \boldsymbol{\delta}_s \boldsymbol{\delta}'_s \right) \mathbf{U}_t \leq \|\boldsymbol{\Delta}\|_{\max}^2 (\max_{t \leq T} \|\mathbf{U}_t\|_1/n)^2$$

$$\zeta_{st}^* \leq \|\mathbf{U}_s\|_{\infty} \|\boldsymbol{\delta}_t\|_{\infty} + \|\mathbf{U}_t\|_{\infty} \|\boldsymbol{\delta}_s\|_{\infty} + \|\boldsymbol{\delta}_t\|_{\infty} \|\boldsymbol{\delta}_s\|_{\infty} \leq 2\|\mathbf{U}\|_{\max} \|\boldsymbol{\Delta}\|_{\max} + \|\boldsymbol{\Delta}\|_{\max}^2 \quad \square$$

Lemma B.4. *Under Assumption 3:*

$$(a) \frac{1}{T} \sum_{t=1}^T \left[\frac{1}{nT} \sum_{s=1}^T \hat{f}_{js} \mathbb{E}(\mathbf{U}'_s \mathbf{U}_t) \right]^2 = O_P(1/T)$$

$$(b) \frac{1}{T} \sum_{t=1}^T \left[\frac{1}{T} \sum_{s=1}^T \hat{f}_{js} \tilde{\zeta}_{st} \right]^2 = O_P[(1/\sqrt{n} + \omega + \omega^2)^2]$$

$$(c) \frac{1}{T} \sum_{t=1}^T \left[\frac{1}{T} \sum_{s=1}^T \hat{f}_{js} \tilde{\eta}_{st} \right]^2 = O_P[(1/\sqrt{n} + \omega)^2]$$

$$(d) \frac{1}{T} \sum_{t=1}^T \left[\frac{1}{T} \sum_{s=1}^T \hat{f}_{js} \xi_{st} \right]^2 = O_P[(1/\sqrt{n} + \omega)^2]$$

$$(e) \frac{1}{T} \sum_{t=1}^T \|\hat{\mathbf{f}}_t - \mathbf{H} \mathbf{f}_t\|^2 = O_P[1/T + (1/\sqrt{n} + \omega + \omega^2)^2]$$

Proof. Part (a) is unaltered by the presence of a pre-estimation so it follows directly from Lemma 8(a) in FLM. For part (b), we have that for $s, l \in [n]$ and $j \in [r]$ by Cauchy-Schwartz inequality

$$\frac{1}{T} \sum_{t=1}^T \left[\frac{1}{T} \sum_{s=1}^T \hat{f}_{js} \tilde{\zeta}_{st} \right]^2 \leq \left[\frac{1}{T^2} \sum_{s,l=1}^T \left(\frac{1}{T} \sum_{t=1}^T \tilde{\zeta}_{st} \tilde{\zeta}_{lt} \right)^2 \right]^{1/2}$$

Since $\tilde{\zeta}_{st} = \zeta_{st} + \zeta_{st}^* = O_P(1/\sqrt{n} + \omega + \omega^2)$ by Lemma B.3, the term in parentheses is $O_P[(1/\sqrt{n} + \omega + \omega^2)^2]$. The result (b) then follows. For (c), by the triangle inequality and Lemma 8(c) in FLM, we have that $\|\sum_{i=1}^n \lambda_{ji} \tilde{u}_{it}\| \leq \|\sum_{i=1}^n \lambda_{ji} U_{it}\| + \|\sum_{i=1}^n \lambda_{ji} \delta_{it}\| = O_P(\sqrt{n}) + O_P(n\omega)$, then we conclude

$$\frac{1}{T} \sum_{t=1}^T \left[\frac{1}{T} \sum_{s=1}^T \hat{\mathbf{f}}_s \tilde{\eta}_{st} \right]^2 \leq \frac{1}{Tn^2} \sum_{t=1}^T \left\| \sum_{i=1}^n \mathbf{U}_{it} \lambda_i \right\|^2 = O_P(1/n + \omega/\sqrt{n} + \omega^2).$$

The proof of part (d) is analogous to part (c) therefore is omitted. For (e), let $[\hat{\mathbf{f}}_t - \mathbf{H} \mathbf{f}_t]_j$ denote the j -th entry of the vector $\hat{\mathbf{f}}_t - \mathbf{H} \mathbf{f}_t$. Since \mathbf{V}/n is bounded away from zero by Lemma B.2(c), the fact that $(a+b+c+d)^2 \leq 4(a^2+b^2+c^2+d^2)$ and using (A.1) we have that $\max_{j \leq r} T^{-1} \sum_t [\hat{\mathbf{f}}_t - \mathbf{H} \mathbf{f}_t]_j$

is upper bounded by some constant $C < \infty$ times

$$\left[\max_{j \leq r} \frac{1}{T} \sum_{t=1}^T \left(\frac{1}{T} \sum_{s=1}^T \hat{f}_{js} \frac{\mathbb{E}(\mathbf{U}'_s \mathbf{U}_t)}{n} \right)^2 + \max_{j \leq r} \frac{1}{T} \sum_{t=1}^T \left(\frac{1}{T} \sum_{s=1}^T \hat{f}_{js} \tilde{\zeta}_{st} \right)^2 \right. \\ \left. + \max_{j \leq r} \frac{1}{T} \sum_{t=1}^T \left(\frac{1}{T} \sum_{s=1}^T \hat{f}_{js} \tilde{\eta}_{st} \right)^2 + \max_{j \leq r} \frac{1}{T} \sum_{t=1}^T \left(\frac{1}{T} \sum_{s=1}^T \hat{f}_{js} \tilde{\xi}_{st} \right)^2 \right].$$

The result then follows by applying the bounds from part (a)-(d) to each of the four terms above. \square

Lemma B.5. *Under Assumption 2:*

- (a) $\max_{t \leq T} \left\| \frac{1}{nT} \sum_{s=1}^T \hat{\mathbf{f}}_s \mathbb{E}(\mathbf{U}'_s \mathbf{U}_t) \right\| = O_P(1/\sqrt{T})$
- (b) $\max_{t \leq T} \left\| \frac{1}{T} \sum_{s=1}^T \hat{\mathbf{f}}_s \tilde{\zeta}_{st} \right\| = O_P(\sqrt{\psi_{p/2}^{-1}(T)/n} + \psi^{-1}(nT)\omega + \omega^2)$
- (c) $\max_{t \leq T} \left\| \frac{1}{T} \sum_{s=1}^T \hat{\mathbf{f}}_s \tilde{\eta}_{st} \right\| = O_P(\psi^{-1}(T)/\sqrt{n} + \omega)$
- (d) $\max_{t \leq T} \left\| \frac{1}{T} \sum_{s=1}^T \hat{\mathbf{f}}_s \tilde{\xi}_{st} \right\| = O_P(\psi^{-1}(T)(1/\sqrt{n} + \omega))$

Proof. Once again, part (a) is unaltered by the presence of a pre-estimation so it follows directly from Lemma 9(a) in FLM. For part (b), from the Cauchy-Schwartz inequality we have

$$\max_{t \leq T} \left\| \frac{1}{T} \sum_{s=1}^T \hat{\mathbf{f}}_s \tilde{\zeta}_{st} \right\| \leq \left(\frac{1}{T} \sum_{s=1}^T \|\hat{\mathbf{f}}_s\|^2 \max_{t \leq T} \frac{1}{T} \sum_{s=1}^T \tilde{\zeta}_{st}^2 \right)^{1/2}.$$

The first summation inside the parentheses equal r due to the normalization. For the second summation, by the triangle inequality, we have $\max_{t \leq T} \frac{1}{T} \sum_{s=1}^T \tilde{\zeta}_{st}^2 \leq \max_{t \leq T} \frac{1}{T} \sum_{s=1}^T \zeta_{st}^2 + 2 \max_{t \leq T} \frac{1}{T} \sum_{s=1}^T \zeta_{st} \zeta_{st}^* + \max_{t \leq T} \frac{1}{T} \sum_{s=1}^T \zeta_{st}^{*2}$. For the first term, the maximum inequality followed by Assumption 2(e) yields

$$\max_{t \leq T} \frac{1}{T} \sum_{s=1}^T \zeta_{st}^2 = O_P \left[\psi_{p/2}^{-1}(T) \max_{s,t} \|\zeta^2\|_{\psi_{p/2}} \right] = O_P \left[\psi_{p/2}^{-1}(T) \max_{s,t} \|\zeta\|_{\psi}^2 \right] = O_P \left[\frac{\psi_{p/2}^{-1}(T)}{n} \right].$$

The last one is $O_P[(\psi^{-1}(nT)\omega + \omega^2)^2]$ by Lemma B.3(d). Then by Cauchy Schwartz we have that $\max_{t \leq T} \frac{1}{T} \sum_{s=1}^T \tilde{\zeta}_{st}^2 = O_P[(\sqrt{\psi_{p/2}^{-1}(T)/n} + \psi^{-1}(nT)\omega + \omega^2)^2]$ and result (b) follows.

For (c), by the triangle inequality we have that $\max_{t \leq T} \left\| \frac{1}{n} \sum_{i=1}^n \boldsymbol{\lambda}_i \tilde{U}_{it} \right\| \leq \max_{t \leq T} \left\| \frac{1}{n} \sum_{i=1}^n \boldsymbol{\lambda}_i U_{it} \right\| + \max_{t \leq T} \left\| \frac{1}{n} \sum_{i=1}^n \boldsymbol{\lambda}_i \delta_{it} \right\|$. For the first term, the maximum inequality followed by Assumption 2(f) yields

$$\max_{t \leq T} \left\| \frac{1}{n} \boldsymbol{\Lambda}' \mathbf{U}_t \right\| = O_P \left[\frac{\psi^{-1}(T)}{\sqrt{n}} \max_t \left\| \frac{1}{\sqrt{n}} \boldsymbol{\Lambda}' \mathbf{U}_t \right\| \right] = O_P(\psi^{-1}(T)/\sqrt{n}).$$

the second term is upper bounded by $r\|\mathbf{\Lambda}\|_{\max}\|\mathbf{\Delta}\|_{\max} = O_P(\omega)$ by Assumption 2(d). We then obtain the result since

$$\max_{t \leq T} \left\| \frac{1}{T} \sum_{s=1}^T \hat{\mathbf{f}}_s \tilde{\eta}_{st} \right\| \leq \left\| \frac{1}{T} \sum_{s=1}^T \hat{\mathbf{f}}_s \mathbf{f}'_s \right\| \max_{t \leq T} \left\| \frac{1}{n} \sum_{i=1}^n \boldsymbol{\lambda}_i \tilde{U}_{it} \right\| = O_P \left(\frac{\psi^{-1}(T)}{\sqrt{n}} + \omega \right). \quad (\text{B.1})$$

By the triangle inequality, $\left\| \frac{1}{nT} \sum_s \sum_i \boldsymbol{\lambda}_i \tilde{U}_{is} \hat{\mathbf{f}}_s \right\| \leq \left\| \frac{1}{nT} \sum_s \sum_i \boldsymbol{\lambda}_i U_{is} \hat{\mathbf{f}}_s \right\| + \left\| \frac{1}{nT} \sum_s \sum_i \boldsymbol{\lambda}_i \delta_{is} \hat{f}_{is} \right\|$. Lemma 9(d) of FLM shows that the first term is $O_P(1/\sqrt{n})$. For the second term for each $j \in [r]$:

$$\left\| \frac{1}{nT} \sum_s \sum_i \boldsymbol{\lambda}_i \delta_{is} \hat{f}_{js} \right\|^2 \leq \left(\frac{1}{T} \sum_{s=1}^n \left\| \frac{1}{n} \sum_{i=1}^n \boldsymbol{\lambda}_i \delta_{is} \right\|^2 \hat{f}_{js} \right) \left(\frac{1}{T} \sum_{s=1}^n \hat{f}_{js}^2 \right) = O_P(\omega^2).$$

Thus $\left\| \frac{1}{nT} \sum_s \sum_i \boldsymbol{\lambda}_i \tilde{U}_{it} \hat{\mathbf{f}}_s \right\| = O_P(1/\sqrt{n} + \omega)$ and by Cauchy-Schwartz inequality we have

$$\max_{t \leq T} \left\| \frac{1}{T} \sum_{s=1}^T \hat{\mathbf{f}}_s \xi_{st} \right\| \leq \max_{t \leq T} \|\mathbf{F}_t\| \left\| \frac{1}{nT} \sum_s \sum_i \boldsymbol{\lambda}_i \tilde{U}_{it} \hat{\mathbf{f}}_s \right\| = O_P(\psi^{-1}(T)(1/\sqrt{n} + \omega)). \quad (\text{B.2})$$

□

Lemma B.6. *Let $\omega_1 + \psi^{-1}(n^2)/\sqrt{T} = O(1)$ where ω_1 is defined in Lemma B.2, then Under Assumption 3 we have*

$$(a) \quad \|\mathbf{V}^{-1}\| = O_P(1/n)$$

$$(b) \quad \|\mathbf{H}\| = O_P(1)$$

$$(c) \quad \|\mathbf{H}'\mathbf{H} - \mathbf{I}_r\|_F = O_P(1/\sqrt{T} + 1/\sqrt{n} + \omega)$$

$$(d) \quad \max_{i \leq n} \frac{1}{T} \sum_{t=1}^T \hat{R}_{it}^2 = O_P(\omega(\psi^{-1}(nT) + \omega) + \psi_{p/2}^{-1}(n)/\sqrt{T} + 1)$$

$$(e) \quad \max_{i \leq n} \max_{j \leq r} \frac{1}{T} \sum_{t=1}^T F_{jt} \tilde{U}_{it} = O_P(\psi_{p/2}^{-1}(n)/\sqrt{T} + \omega)$$

Proof. We have that $\mathbf{V}^{-1} = \text{diag}(1/\hat{a}_1, \dots, 1/\hat{a}_r)$ and $1/\hat{a}_j \asymp_P 1/n$ for $j \leq r$ by Lemma B.2(c). The result (a) then follows. The normalization tell us $\|\hat{\mathbf{F}}\| = \sqrt{T}$, Lemma 11(a) in FLM give us $\|\mathbf{F}\| = O_P(\sqrt{T})$, $\|\mathbf{\Lambda}'\mathbf{\Lambda}\| = \tilde{a}_1 \asymp n$ by Lemma B.1(a) and from part (b) we have $\|\mathbf{V}^{-1}\| = O_P(1/n)$. Result (b) then follows since by definition $\mathbf{H} := T^{-1}\mathbf{V}^{-1}\hat{\mathbf{F}}'\mathbf{F}\mathbf{\Lambda}'\mathbf{\Lambda}$. For (c) we have by the triangle inequality

$$\|\mathbf{H}'\mathbf{H} - \mathbf{I}_r\|_F \leq \|\mathbf{H}'\mathbf{H} - \mathbf{H}'\mathbf{F}'\mathbf{F}/T\mathbf{H}\|_F + \|\mathbf{H}'\mathbf{F}'\mathbf{F}/T\mathbf{H} - \mathbf{I}_r\|_F$$

For the first term we have

$$\|\mathbf{H}'(\mathbf{I}_r - \mathbf{F}'\mathbf{F}/T)\mathbf{H}\|_F \leq \|\mathbf{H}\|^2 \|\mathbf{I}_r - \mathbf{F}'\mathbf{F}/T\|_F = O_P(1)O_P(1/\sqrt{T}).$$

The second term is equal to

$$\|\mathbf{H}'\mathbf{F}'\mathbf{F}/T\mathbf{H} - \widehat{\mathbf{F}}'\widehat{\mathbf{F}}/T\|_F$$

For (d) we have

$$\begin{aligned} \max_{i \leq n} \frac{1}{T} \sum_{t=1}^T \widehat{R}_{it}^2 &\leq \max_{i \leq n} \frac{1}{T} \sum_{t=1}^T (\widehat{R}_{it}^2 - R_{it}^2) + \max_{i \leq n} \frac{1}{T} \sum_{t=1}^T R_{it}^2 - \mathbb{E}(R_{it}^2) + \max_{i \leq n} \frac{1}{T} \sum_{t=1}^T \mathbb{E}(R_{it}^2) \\ &\leq \max_{i,t} |\widehat{R}_{it}^2 - R_{it}^2| + \max_{i \leq n} \frac{1}{T} \sum_{t=1}^T R_{it}^2 - \mathbb{E}(R_{it}^2) + \max_{i,t} \mathbb{E}(R_{it}^2). \end{aligned}$$

The last term is $O(1)$ by Assumption 3(a), the middle term $O_P(\psi_{p/2}^{-1}(n)/\sqrt{T})$. The first term is no larger than $\|\Delta\|_{\max}(2\|\mathbf{R}\|_{\max} + \|\Delta\|_{\max}) = O_P(\omega(\psi^{-1}(nT) + \omega))$. The result (d) then follows.

For (e) we have for each $j \leq r$:

$$\begin{aligned} |T^{-1} \sum_t F_{jt} \widetilde{U}_{it}| &\leq |T^{-1} \sum_t F_{jt} U_{it}| + |T^{-1} \sum_t F_{jt} \delta_{it}| \\ &\leq |T^{-1} \sum_t F_{jt} U_{it}| + (T^{-1} \sum_t F_{jt}^2 T^{-1} \sum_t \delta_{it}^2)^{1/2} \end{aligned}$$

The first term is $O_P(\psi_{p/2}^{-1}(n)/\sqrt{T})$ by the maximum inequality and Assumption 3 and the second is $O_P(\omega)$. \square

Lemma B.7. For every $s > 0$:

$$\rho_p(S, Z) \leq \rho_p(\widetilde{T}, \widetilde{Z}) + \Delta_p(\sqrt{\frac{mq}{n}} \widetilde{Z}, s) + \rho_p(\sqrt{\frac{mq}{n}} \widetilde{Z}, Z) + \mathbb{P}(\sqrt{\frac{mr}{n}} \|\widetilde{U}\|_p > s) + \rho_p(T, \widetilde{T}) + \rho_p(U, \widetilde{U}).$$

Proof. We start by showing that for every pair of random variables X and Y defined in the same probability space taking values in the normed space $(S, \|\cdot\|)$ and pair of non-negative reals t, s , we have

$$\mathbb{P}(\|X\| \leq t - s) - \mathbb{P}(\|Y\| > s) \leq \mathbb{P}(\|X + Y\| \leq t) \leq \mathbb{P}(\|X\| \leq t + s) + \mathbb{P}(\|Y\| > s). \quad (\text{B.3})$$

Indeed, for the right hand side inequality we use $\|X + Y\| = \|X - (-Y)\| \geq \|X\| - \|Y\|$. Hence, for

any $t, s > 0$:

$$\begin{aligned}
\mathbb{P}(\|X + Y\| \leq t) &\leq \mathbb{P}(\|X\| \leq t + \|Y\|) \\
&\leq \mathbb{P}(\|X\| \leq t + \|Y\|, \|Y\| \leq s) + \mathbb{P}(\|Y\| > s) \\
&\leq \mathbb{P}(\|X\| \leq t + s) + \mathbb{P}(\|Y\| > s).
\end{aligned}$$

For the other side we use $\|X + Y\| \leq \|X\| + \|Y\|$ to write

$$\begin{aligned}
\mathbb{P}(\|X + Y\| \leq t) &\geq \mathbb{P}(\|X\| \leq t - \|Y\|) \\
&\geq \mathbb{P}(\|X\| \leq t - \|Y\|) + \mathbb{P}(\|Y\| > s) - \mathbb{P}(\|Y\| > s)
\end{aligned}$$

Now replace X and Y by $\sqrt{\frac{mq}{n}}T$ and $\sqrt{\frac{mr}{n}}U$ in (B.3), respectively and set $\|\cdot\| = \|\cdot\|_p$. The right hand side of the resulting expression can be upper bounded by $\mathbb{P}(\sqrt{\frac{mq}{n}}\|\tilde{T}\|_p \leq t + s) + \mathbb{P}(\sqrt{\frac{mr}{n}}\|\tilde{U}\| > s) + \rho_p(T, \tilde{T}) + \rho_p(U, \tilde{U})$, whereas the left hand side can be lower bounded by $\mathbb{P}(\sqrt{\frac{mq}{n}}\|\tilde{T}\| \leq t - s) - \mathbb{P}(\sqrt{\frac{mr}{n}}\|\tilde{U}\| > s) - \rho_p(T, \tilde{T}) - \rho_p(U, \tilde{U})$. Therefore

$$\begin{aligned}
&\mathbb{P}(\sqrt{\frac{mq}{n}}\|\tilde{T}\|_p \leq t - s) - \mathbb{P}(\sqrt{\frac{mr}{n}}\|\tilde{U}\|_p > s) - \rho_p(T, \tilde{T}) - \rho_p(U, \tilde{U}) \\
&\leq \mathbb{P}(\|S\|_p \leq t) \\
&\leq \mathbb{P}(\sqrt{\frac{mq}{n}}\|\tilde{T}\|_p \leq t + s) + \mathbb{P}(\sqrt{\frac{mr}{n}}\|\tilde{U}\|_p > s) + \rho_p(T, \tilde{T}) + \rho_p(U, \tilde{U}).
\end{aligned}$$

Then for the right-hand side

$$\begin{aligned}
\mathbb{P}(\sqrt{\frac{mq}{n}}\|\tilde{T}\|_p \leq t + s) &\leq \mathbb{P}(\sqrt{\frac{mq}{n}}\|\tilde{Z}\|_p \leq t + s) + \rho_p(\tilde{T}, \tilde{Z}) \\
&\leq \mathbb{P}(\sqrt{\frac{mq}{n}}\|\tilde{Z}\|_p \leq t) + \Delta_p(\sqrt{\frac{mq}{n}}\tilde{Z}, s) + \rho_p(\tilde{T}, \tilde{Z}) \\
&\leq \mathbb{P}(\|Z\|_p \leq t) + \rho_p(\sqrt{\frac{mq}{n}}\tilde{Z}, Z) + \Delta_p(\sqrt{\frac{mq}{n}}\tilde{Z}, s) + \rho_p(\tilde{T}, \tilde{Z})
\end{aligned}$$

Similarly for the left-hand side and the proof is completed. □

By the triangle inequality $\|\hat{\mathbf{Y}} - \mathbf{Y}\|_{\max} \leq \|\hat{\mathbf{Y}} - \tilde{\mathbf{Y}}\|_{\max} + \|\tilde{\mathbf{Y}} - \mathbf{Y}\|_{\max}$ where $\tilde{\mathbf{Y}}$ is the sample

covariance matrix of $\tilde{\mathbf{D}}_t := U_{1t}\mathbf{U}_{-1t}$. The second term is $O(\psi_{p/4}^{-1}(n^2)/\sqrt{T})$ while for the first

$$\|\hat{\mathbf{Y}} - \tilde{\mathbf{Y}}\|_{\max} \leq \|\mathbf{D} - \tilde{\mathbf{D}}\|_{\max}(2\|\tilde{\mathbf{D}}\|_{\max} + \|\mathbf{D} - \tilde{\mathbf{D}}\|_{\max})$$

The first term in parentheses is $O(\psi_*^{-1}(nT))$ and the second can be upper bounded by $\|\hat{\mathbf{U}} - \mathbf{U}\|_{\max}(2\|\mathbf{U}\|_{\max} + \|\hat{\mathbf{U}} - \mathbf{U}\|_{\max})$ which is show to be $O_P(\eta(n, T)\psi^{-1}(nT))$ in the proof of Lemma B.16. Therefore we conclude

$$\|\hat{\mathbf{Y}} - \mathbf{Y}\|_{\max} = O_P\left(\eta(n, T)\psi^{-1}(nT)\psi_{p/2}^{-1}(nT) + \psi_{p/4}^{-1}(n^2)/\sqrt{T}\right)$$

To leverage on the results of Gaussian approximation, in particular on the work of Giessing and Fan (2020) we would like to establish some sort of asymptotic linearity namely

$$\mathbf{Q}_T = \frac{1}{\sqrt{T}} \sum_{t=1}^T \mathbf{D}_t = \frac{1}{\sqrt{T}} \sum_{t=1}^T \tilde{\mathbf{D}}_t + \mathbf{R}_T =: \tilde{\mathbf{Q}}_T + \mathbf{R}_T. \quad (\text{B.4})$$

such that $\|\mathbf{R}_t\|_{\infty}$ vanishes in probability at an appropriate rate as $n, T \rightarrow \infty$. Then we can approximate the distribution of $S = \|\mathbf{Q}\|_{\infty}$ by the distribution of $\tilde{S} := \|\tilde{\mathbf{Q}}\|_p$, which in turn can be approximated by the distribution of $S^* := \|\mathbf{Q}^*\|_{\infty}$ with high probability.

For some $\epsilon > 0$ we might set

$$\begin{aligned} \delta_1 &= h[\eta(n, T)(\psi^{-1}(nT))^3 + \psi_{p/4}^{-1}(n^2)/\sqrt{T}] \\ \delta_2 &= \eta^{1-\epsilon}[\psi^{-1}(n) + \sqrt{T}\eta] \end{aligned}$$

Lemma B.8. $\|\hat{\mathbf{Y}} - \mathbf{Y}\|_{\max} = O_P\left(h[\eta(\psi_p^{-1}(nT))^3 + \psi_{p/4}^{-1}(n^4)/\sqrt{T}]\right)$

Proof. Let $\mathbf{i} := (i_1, i_2, i_3, i_4)$ be a multi-index where $i_1, i_2, i_3, i_4 \in [n]$. Define for \mathbf{i} and $|\ell| < T$:

$$\tilde{\gamma}_{\mathbf{i}}^{\ell} := \frac{1}{T} \sum_{t=|\ell|+1}^T U_{i_1,t} U_{i_2,t} U_{i_3,t-|\ell|} U_{i_4,t-|\ell|}; \quad \gamma_{\mathbf{i}}^{\ell} := \mathbb{E}\tilde{\gamma}_{\mathbf{i}}^{\ell},$$

and $\hat{\gamma}_{\mathbf{i}}^{\ell}$ as $\tilde{\gamma}_{\mathbf{i}}^{\ell}$ with U 's replaced by \hat{U} 's. Also define

$$\tilde{v}_{\mathbf{i}} := \sum_{|\ell| < T} k(\ell/h) \tilde{\gamma}_{\mathbf{i}}^{\ell} \quad v_{\mathbf{i}} := \sum_{|\ell| < T} \gamma_{\mathbf{i}}^{\ell},$$

and \hat{v}_i as \tilde{v}_i with U 's replaced by \hat{U} 's. Then we write

$$\tilde{v}_i - v_i = \sum_{|\ell| < T} k(\ell/h)(\tilde{\gamma}_i^\ell - \gamma_i^\ell) + \sum_{|\ell| < T} (k(\ell/h) - 1)\gamma_i^\ell. \quad (\text{B.5})$$

Since $\|\tilde{\gamma}_i^\ell - \gamma_i^\ell\|_{\psi_{p/4}} = O(\sqrt{T - |\ell|}/T) = O(1/\sqrt{T})$, the $\psi_{p/4}$ -Orlicz norm of the first term is bounded by

$$h \sum_{|\ell| < T} |h^{-1}k(\ell/h)| \|\tilde{\gamma}_i^\ell - \gamma_i^\ell\|_{\psi_{p/4}} = O\left(\frac{h}{\sqrt{T}} \int |k(u)| du\right) = O(h/\sqrt{T}),$$

whereas the second term is deterministic and is shown to be $O(h/\sqrt{T})$ by Andrews (1991). Thus $\|\tilde{v}_i - v_i\|_{\psi_{p/4}} = O(h/\sqrt{T})$ uniformly in $i \in [n]^4$. Thus, by the maximal inequality followed by Markov's inequality we conclude that

$$\max_i |\tilde{v}_i - v_i| = O_P(\psi_{p/4}^{-1}(n^4) \max_i \|\tilde{v}_i - v_i\|_{\psi_{p/4}}) = O_P[\psi_{p/4}^{-1}(n^4)h/\sqrt{T}]. \quad (\text{B.6})$$

We now use the fact that for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^q$ we have $|\prod_{i=1}^q x_i - \prod_{i=1}^q y_i| = O(\sum_{i=0}^{q-1} \|\mathbf{x} - \mathbf{y}\|_\infty^{n-i} \|\mathbf{y}\|_\infty^i)$ combined with the fact that $\|\hat{\mathbf{U}} - \mathbf{U}\|_{\max} = o(1)$ to obtain

$$\begin{aligned} \max_{i,\ell} |\hat{\gamma}_i^\ell - \tilde{\gamma}_i^\ell| &\leq \max_{i,t,\ell} |\hat{U}_{i_1,t} \hat{U}_{i_2,t} \hat{U}_{i_3,t-|\ell|} \hat{U}_{i_4,t-|\ell|} - U_{i_1,t} U_{i_2,t} U_{i_3,t-|\ell|} U_{i_4,t-|\ell|}| \\ &= O(\|\hat{\mathbf{U}} - \mathbf{U}\|_{\max} \|\mathbf{U}\|_{\max}^3) \\ &= O_P[\eta[\psi_p^{-1}(nT)]^3] \end{aligned}$$

Therefore we conclude

$$\max_i |\hat{v}_i - \tilde{v}_i| \leq \max_{i,\ell} |\hat{\gamma}_i^\ell - \tilde{\gamma}_i^\ell| \sum_{|\ell| < T} |k(\ell/h)| = O_P\left(h\eta[\psi_p^{-1}(nT)]^3 \int |k(u)| du\right) = O_P(h\eta[\psi_p^{-1}(nT)]^3). \quad (\text{B.7})$$

The result then follows from the triangle inequality $\|\hat{\mathbf{Y}} - \mathbf{Y}\|_{\max} \leq \max_i |\hat{v}_i - \tilde{v}_i| + \max_i |\tilde{v}_i - v_i|$, expression (B.10) and (B.11). \square

Lemma B.9. *If $\|\delta_{it}\|_{\psi_p} \leq C < \infty$ where $\delta_{it} := \hat{R}_{it} - R_{it}$ then*

$$\|(\mathbf{V}/n)(\mathbf{F}_t - \mathbf{H}\mathbf{F}_t)\|_2\|_{\psi_p} = O\left(\frac{1}{\sqrt{T}} + \frac{\psi_{p/2}^{-1}(T)}{\sqrt{n}} + \psi_{p/2}^{-1}(T)C\right).$$

Proof. In this proof we use the fact that for any (possibly random) A_{st} , by Cauchy-Schwartz inequality and the normalization $\widehat{\mathbf{F}}\widehat{\mathbf{F}}/T = \mathbf{I}_r$, we have $\|\frac{1}{T}\sum_{s=1}^T \widehat{\mathbf{F}}_s A_{st}\| \leq \sqrt{r} \left(\frac{1}{T}\sum_{s=1}^T A_{st}^2\right)^{1/2}$. Thus

$$g(A_{st}) := \left\| \frac{1}{T} \sum_{s=1}^T \widehat{\mathbf{F}}_s A_{st} \right\|_{\psi} = O \left[\left\| \left(\frac{1}{T} \sum_{s=1}^T A_{st}^2 \right)^{1/2} \right\|_{\psi} \right].$$

(a) Set $A_{st} = \mathbb{E}(\mathbf{U}'_s \mathbf{U}_t)/n$, then $g(A_{st}) = O(1/\sqrt{T})$.

(b) Set $A_{st} = \tilde{\zeta}_{st} := (\mathbf{U}'_s \mathbf{U}_t - \mathbb{E}(\mathbf{U}'_s \mathbf{U}_t))/n$, then by maximal inequality $g(A_{st}) = O(\|\max_{s \leq T} \tilde{\zeta}_{st}\|_{\psi}) = O(\psi^{-1}(T) \max_{s \leq T} \|\tilde{\zeta}_{st}\|_{\psi})$. By the triangle inequality $\|\tilde{\zeta}_{st}\|_{\psi} \leq \|\zeta_{st}\|_{\psi} + \|\zeta_{st}^*\|_{\psi}$. The first term is $O(1/\sqrt{n})$ by Assumption 3(d). The second can be upper bounded by $\|\mathbf{U}'_s \boldsymbol{\delta}_t/n\|_{\psi} + \|\boldsymbol{\delta}'_s \mathbf{U}_t/n\|_{\psi} + \|\boldsymbol{\delta}'_s \boldsymbol{\delta}_t/n\|_{\psi} = O(\|U_{is}\|_{\psi_{p/2}} \|\delta_{it}\|_{\psi_{p/2}}) + O(\|\delta_{it}\|_{\psi_{p/2}}^2)$. Thus $g(A_{st}) = O(\psi^{-1}(T)(1/\sqrt{n} + C + C^2))$.

(c) Set $A_{st} = \tilde{\eta}_{st} := \mathbf{F}'_s \sum_{i=1}^n \boldsymbol{\lambda}_i (U_{it} + \delta_{it})/n$, then apply Cauchy-Schwartz twice to obtain

$$g(A_{st}) = O\left(\left(\frac{1}{T} \sum_{s=1}^T \|\mathbf{F}_s\|^2\right)^{1/2} \left\| \sum_{i=1}^n \boldsymbol{\lambda}_i \frac{U_{it} + \delta_{it}}{n} \right\|_{\psi_{p/2}}\right) = O(1)O\left(\left\| \sum_{i=1}^n \boldsymbol{\lambda}_i \frac{U_{it}}{n} \right\|_{\psi_{p/2}} + \left\| \sum_{i=1}^n \boldsymbol{\lambda}_i \frac{\delta_{it}}{n} \right\|_{\psi_{p/2}}\right).$$

The first term in square brackets is $O(1/\sqrt{n})$ by Assumption 2(d) and 3(e); the second is $O(C)$. Hence $g(A_{st}) = O(\frac{1}{\sqrt{n}} + C)$.

(d) Set $A_{st} = \tilde{\xi}_{st} := \mathbf{F}'_t \sum_{i=1}^n \boldsymbol{\lambda}_i (U_{is} + \delta_{is})/n$, then apply Cauchy-Schwartz twice followed by the maximal inequality to obtain

$$\begin{aligned} g(A_{st}) &= O\left(\|\mathbf{F}_t\|_{\psi_{p/2}} \left\| \left(\frac{1}{T} \sum_{s=1}^T \left\| \sum_{i=1}^n \boldsymbol{\lambda}_i \frac{U_{is} + \delta_{is}}{n} \right\|^2 \right)^{1/2} \right\|_{\psi_{p/2}}\right) \\ &= O(1)O(\psi^{-1}(T) \left[\left\| \sum_{i=1}^n \boldsymbol{\lambda}_i \frac{U_{is}}{n} \right\|_{\psi_{p/2}} + \left\| \sum_{i=1}^n \boldsymbol{\lambda}_i \frac{\delta_{is}}{n} \right\|_{\psi_{p/2}} \right]). \end{aligned}$$

The first term in square brackets is $O(1/\sqrt{n})$ by Assumption 2(d) and 3(e); the second is $O(C)$. Hence $g(A_{st}) = O(\psi_{p/2}^{-1}(T)[\frac{1}{\sqrt{n}} + C])$.

Finally, use the identity (A.1), the triangle inequality twice and the bounds (a) – (d) to obtain the result. \square

Lemma B.10. *If $\max_{i,t} \|\delta_{it}\|_{\psi} = O(C)$ and $\|\widehat{\mathbf{U}} - \mathbf{U}\|_{\max} = O_P(\eta)$ then*

$$\left\| \frac{1}{\sqrt{T}} (\widehat{\mathbf{U}}\widehat{\mathbf{U}}' - \mathbf{U}\mathbf{U}') \right\|_{\max} = O_P \left(\sqrt{T}\eta^2 + \frac{r_1}{\sqrt{T}} + \frac{r_2}{\sqrt{n}} + r_3 C \right)$$

where

$$\begin{aligned} r_1 &:= \psi_p^{-1}(n)\psi_{p/2}^{-1}(n)\psi_{p/2}^{-1}(n^2) \\ r_2 &:= \psi_p^{-1}(n)\psi_{p/4}^{-1}(T) + \psi_{p/2}^{-1}(n) \\ r_3 &:= \psi_p^{-1}(n)\psi_{p/4}^{-1}(T) + \psi_p^{-1}(nT)\psi_{p/2}^{-1}(n). \end{aligned}$$

Proof. By the triangle inequality we have

$$\left\| \frac{1}{\sqrt{T}}(\hat{\mathbf{U}}\hat{\mathbf{U}}' - \mathbf{U}\mathbf{U}') \right\|_{\max} \leq \left\| \frac{1}{\sqrt{T}}(\hat{\mathbf{U}} - \mathbf{U})(\hat{\mathbf{U}} - \mathbf{U})' \right\|_{\max} + 2 \left\| \frac{1}{\sqrt{T}}\mathbf{U}(\hat{\mathbf{U}} - \mathbf{U})' \right\|_{\max}.$$

For the first term we have

$$\left\| \frac{1}{\sqrt{T}}(\hat{\mathbf{U}} - \mathbf{U})(\hat{\mathbf{U}} - \mathbf{U})' \right\|_{\max} \leq \sqrt{T}\|\hat{\mathbf{U}} - \mathbf{U}\|_{\max}^2 = O_P(\sqrt{T}\eta^2).$$

For the second term we use decomposition (A.3) to write

$$\begin{aligned} \frac{1}{\sqrt{T}} \sum_{t=1}^T U_{it}(\hat{U}_{jt} - U_{jt}) &= \frac{1}{\sqrt{T}} \sum_{t=1}^T U_{it}(\hat{\boldsymbol{\lambda}}_j' \hat{\mathbf{F}}_t - \boldsymbol{\lambda}_j' \mathbf{F}_t + \hat{R}_{jt} - R_{jt}) \\ &= \left[(\hat{\boldsymbol{\lambda}}_j - \mathbf{H}\boldsymbol{\lambda}_j) + \mathbf{H}\boldsymbol{\lambda}_j \right]' \frac{1}{\sqrt{T}} \sum_{t=1}^T U_{it}(\hat{\mathbf{F}}_t - \mathbf{H}\mathbf{F}_t) \\ &\quad + \left[(\hat{\boldsymbol{\lambda}}_j - \mathbf{H}\boldsymbol{\lambda}_j) + (\mathbf{H}'\mathbf{H} - I_r)\boldsymbol{\lambda}_j \right]' \frac{1}{\sqrt{T}} \sum_{t=1}^T U_{it}\mathbf{F}_t + (\hat{\gamma}_j - \gamma_j)' \frac{1}{\sqrt{T}} \sum_{t=1}^T U_{it}\mathbf{W}_{jt} \end{aligned}$$

Apply Cauchy-Schwartz inequality in each term followed by the triangle inequality we obtain

$$\begin{aligned} \left\| \frac{1}{\sqrt{T}}\mathbf{U}(\hat{\mathbf{U}} - \mathbf{U})' \right\|_{\max} &\leq \left[\max_{j \leq n} \|\hat{\boldsymbol{\lambda}}_j - \mathbf{H}\boldsymbol{\lambda}_j\| + \sqrt{r}\|\mathbf{H}\|\|\boldsymbol{\Lambda}\|_{\max} \right] \max_{i \leq n} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T U_{it}(\hat{\mathbf{F}}_t - \mathbf{H}\mathbf{F}_t) \right\| \\ &\quad + \left[\max_{j \leq n} \|\hat{\boldsymbol{\lambda}}_j - \mathbf{H}\boldsymbol{\lambda}_j\| + \sqrt{r}\|\mathbf{H}'\mathbf{H} - I_r\|\|\boldsymbol{\Lambda}\|_{\max} \right] \max_{i \leq n} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T U_{it}\mathbf{F}_t \right\| \\ &\quad + \max_{j \leq n} \|\hat{\gamma}_j - \gamma_j\| \max_{i, j \leq n} \left\| \frac{1}{\sqrt{T}} \sum_{t=1}^T U_{it}\mathbf{W}_{jt} \right\|. \end{aligned}$$

The first term is $O_P(1)O_P(\psi^{-1}(n)[\frac{1}{\sqrt{T}} + \frac{\psi_{p/2}^{-1}(T)}{\sqrt{n}} + \psi_{p/2}^{-1}(T)C])$ due to Lemma B.6(a), Lemma B.9 and the maximal inequality; the second term is $O_P(\frac{\psi_{p/2}^{-1}(n)}{\sqrt{T}} + \frac{1}{\sqrt{n}} + \psi^{-1}(nT)C)O_P(\psi_{p/2}^{-1}(n))$ since, by the maximal inequality, we might take $\omega = \psi^{-1}(nT)C$ in Theorem 2(b). The last term is

$O_P(\psi^{-1}(n)\psi_{p/2}^{-1}(n)/\sqrt{T})O_P(\psi_{p/2}^{-1}(n^2))$. Thus, $\left\|\frac{1}{\sqrt{T}}\mathbf{U}(\widehat{\mathbf{U}} - \mathbf{U})'\right\|_{\max} = O_P(r_4)$ where

$$r_4 := \frac{\psi_p^{-1}(n)\psi_{p/2}^{-1}(n)\psi_{p/2}^{-1}(n^2)}{\sqrt{T}} + \frac{\psi_p^{-1}(n)\psi_{p/4}^{-1}(T) + \psi_{p/2}^{-1}(n)}{\sqrt{n}} + (\psi_p^{-1}(n)\psi_{p/4}^{-1}(T) + \psi_p^{-1}(nT)\psi_{p/2}^{-1}(n))C. \quad (\text{B.8})$$

The result then follows. \square

Lemma B.11. *Let $\|\widehat{\mathbf{U}} - \mathbf{U}\| = O_P(\eta)$ then $\max_{i,j,t} |\widehat{V}_{ij,t} - V_{ij,t}| = O_P(s_0[\eta + \xi\psi^{-1}(n)])$.*

Proof. By the triangle inequality we have

$$|\widehat{V}_{ij,t} - V_{ij,t}| \leq |\widehat{U}_{i,t} - U_{i,t}| + |\widehat{\boldsymbol{\theta}}_i' \widehat{\mathbf{U}}_{-ij,t} - \boldsymbol{\theta}_i' \mathbf{U}_{-ij,t}|.$$

Using Hölder's inequality, the second term can be further bounded as

$$\begin{aligned} |\widehat{\boldsymbol{\theta}}_i' \widehat{\mathbf{U}}_{-ij,t} - \boldsymbol{\theta}_i' \mathbf{U}_{-ij,t}| &\leq |\widehat{\boldsymbol{\theta}}_i' (\widehat{\mathbf{U}}_{-ij,t} - \mathbf{U}_{-ij,t})| + |(\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i)' \mathbf{U}_{-ij,t}| \\ &\leq \|\widehat{\boldsymbol{\theta}}_i\|_1 \|\widehat{\mathbf{U}}_{-ij,t} - \mathbf{U}_{-ij,t}\|_\infty + \|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i\|_1 \|\mathbf{U}_{-ij,t}\|_\infty \\ &\leq (\|\boldsymbol{\theta}_i\|_1 + \|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i\|_1) \|\widehat{\mathbf{U}}_{-ij,t} - \mathbf{U}_{-ij,t}\|_\infty + \|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i\|_1 \|\mathbf{U}_{-ij,t}\|_\infty. \end{aligned}$$

Combining the last two expressions with the fact that $\|\boldsymbol{\theta}_i\|_1 \leq s_0 \|\boldsymbol{\theta}_i\|_\infty \leq Cs_0$ and $\|\widehat{\boldsymbol{\theta}}_1 - \boldsymbol{\theta}_1\|_1 = O_P(\xi s_0) = O_P(1)$ by Assumption 3(f) and the the maximum inequality yields the result \square

Lemma B.12. *Let $\|\widehat{\mathbf{U}} - \mathbf{U}\| = O_P(\eta)$ then*

$$\max_{i,j} \left| \frac{1}{\sqrt{T}} \sum_{t=1}^T (\widehat{V}_{ij,t} \widehat{V}_{ji,t} - V_{ij,t} V_{ji,t}) \right| = O_P\{s_0^2[r_4 + \xi\psi^{-1}(n) + \sqrt{T}(\eta + \xi\psi^{-1}(n))]^2\}$$

Proof. By the triangle inequality

$$\max_{i,j} \left| \frac{1}{\sqrt{T}} (\widehat{\mathbf{V}}_{ij}' \widehat{\mathbf{V}}_{ji} - \mathbf{V}_{ij}' \mathbf{V}_{ji}) \right| \leq \max_{i,j} \left| \frac{1}{\sqrt{T}} (\widehat{\mathbf{V}}_{ij} - \mathbf{V}_{ij}) (\widehat{\mathbf{V}}_{ji} - \mathbf{V}_{ji}) \right| + 2 \max_{i,j} \left| \frac{1}{\sqrt{T}} \mathbf{V}_{ij}' (\widehat{\mathbf{V}}_{ij} - \mathbf{V}_{ij}) \right|.$$

The first term can be bounded using Lemma B.11 since

$$\max_{i,j} \left| \frac{1}{\sqrt{T}} (\widehat{\mathbf{V}}_{ij} - \mathbf{V}_{ij}) (\widehat{\mathbf{V}}_{ji} - \mathbf{V}_{ji}) \right| \leq \sqrt{T} [\max_{i,j,t} |\widehat{V}_{ij,t} - V_{ij,t}|]^2 = O_P(\sqrt{T}[s_0(\eta + \xi\psi^{-1}(n))]^2).$$

The second term can be upper bounded by

$$\begin{aligned} \max_{i,j} \left| \frac{1}{\sqrt{T}} \mathbf{V}'_{ij} (\hat{\mathbf{U}}_i - \mathbf{U}_i) \right| + \max_{i,j} \|\hat{\boldsymbol{\theta}}_{ij}\|_1 \max_{i,j} \left\| \frac{1}{\sqrt{T}} \mathbf{V}'_{ij} (\hat{\mathbf{U}}_{-ij} - \mathbf{U}_{-ij}) \right\|_\infty \\ + \max_{i,j} \|\hat{\boldsymbol{\theta}}_{ij} - \boldsymbol{\theta}_{ij}\|_1 \max_{i,j} \left| \frac{1}{\sqrt{T}} \mathbf{V}'_{ij} \mathbf{U}_{-ij} \right|. \end{aligned}$$

Recall the rate r_4 appearing in (B.8). Then the first term is $O_P(s_0 r_4)$, the second $O_P(s_0^2 r_4)$ and the last term is $O_P(\xi s_0^2 \psi^{-1}(n))$. Thus $\max_{i,j} \left| \frac{1}{\sqrt{T}} \mathbf{V}'_{ij} (\hat{\mathbf{V}}_{ij} - \mathbf{V}_{ij}) \right| = O_P[s_0^2 (r_4 + \xi \psi^{-1}(n))]$. The result then follows. \square

Lemma B.13. $\|\hat{\boldsymbol{\Upsilon}}_V - \boldsymbol{\Upsilon}_V\|_{\max} = O_P \left(h [s_0 [\eta + \xi \psi_p^{-1}(n)] (s_0 \psi_p^{-1}(nT))^3 + s_0 \frac{\psi_p^{-1}(n^4)}{\sqrt{T}}] \right)$

Proof. The proof is similar to the proof of Lemma B.8, refer to it for details. It suffices to bound in probability $\|\hat{\mathbf{V}} - \mathbf{V}\|_{\max}$ and $\|\mathbf{V}\|_{\max}$, where \mathbf{V} is $(n^2 \times T)$ matrix whose entries are $V_{ij,t}$ for $i, j \in [n]$ and $t \in [T]$. Similar for $\hat{\mathbf{V}}$ with $V_{ij,t}$ replaced $\hat{V}_{ij,t}$. Lemma B.11 bounds the former, for the later we have $\|\mathbf{V}\|_{\max} \leq \max_{i,j} \|\boldsymbol{\theta}_{ij}\|_1 \|\mathbf{U}\|_{\max} = O(s_0 \psi^{-1}(nT))$.

Let $\mathbf{i} := (i_1, i_2, i_3, i_4)$ be a multi-index where $i_1, i_2, i_3, i_4 \in [n]$. Define for \mathbf{i} and $|\ell| < T$:

$$\tilde{\gamma}_{\mathbf{i}}^\ell := \frac{1}{T} \sum_{t=|\ell|+1}^T U_{i_1,t} U_{i_2,t} U_{i_3,t-|\ell|} U_{i_4,t-|\ell|}; \quad \gamma_{\mathbf{i}}^\ell := \mathbb{E} \tilde{\gamma}_{\mathbf{i}}^\ell,$$

and $\hat{\gamma}_{\mathbf{i}}^\ell$ as $\tilde{\gamma}_{\mathbf{i}}^\ell$ with U 's replaced by \hat{U} 's. Also define

$$\tilde{v}_{\mathbf{i}} := \sum_{|\ell| < T} k(\ell/h) \tilde{\gamma}_{\mathbf{i}}^\ell \quad v_{\mathbf{i}} := \sum_{|\ell| < T} \gamma_{\mathbf{i}}^\ell,$$

and $\hat{v}_{\mathbf{i}}$ as $\tilde{v}_{\mathbf{i}}$ with U 's replaced by \hat{U} 's. Then we write

$$\tilde{v}_{\mathbf{i}} - v_{\mathbf{i}} = \sum_{|\ell| < T} k(\ell/h) (\tilde{\gamma}_{\mathbf{i}}^\ell - \gamma_{\mathbf{i}}^\ell) + \sum_{|\ell| < T} (k(\ell/h) - 1) \gamma_{\mathbf{i}}^\ell. \quad (\text{B.9})$$

Since $\|\tilde{\gamma}_{\mathbf{i}}^\ell - \gamma_{\mathbf{i}}^\ell\|_{\psi_{p/4}} = O(\sqrt{T - |\ell|}/T) = O(1/\sqrt{T})$, the $\psi_{p/4}$ -Orlicz norm of the first term is bounded by

$$h \sum_{|\ell| < T} |h^{-1} k(\ell/h)| \|\tilde{\gamma}_{\mathbf{i}}^\ell - \gamma_{\mathbf{i}}^\ell\|_{\psi_{p/4}} = O \left(\frac{h}{\sqrt{T}} \int |k(u)| du \right) = O(h/\sqrt{T}),$$

whereas the second term is deterministic and is shown to be $O(h/\sqrt{T})$ by Andrews (1991). Thus $\|\tilde{v}_{\mathbf{i}} - v_{\mathbf{i}}\|_{\psi_{p/4}} = O(h/\sqrt{T})$ uniformly in $\mathbf{i} \in [n]^4$. Thus, by the maximal inequality followed by

Markov's inequality we conclude that

$$\max_i |\tilde{v}_i - v_i| = O_P(\psi_{p/4}^{-1}(n^4) \max_i \|\tilde{v}_i - v_i\|_{\psi_{p/4}}) = O_P[\psi_{p/4}^{-1}(n^4)h/\sqrt{T}]. \quad (\text{B.10})$$

We now use the fact that for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^q$ we have $|\prod_{i=1}^q x_i - \prod_{i=1}^q y_i| = O(\sum_{i=0}^{q-1} \|\mathbf{x} - \mathbf{y}\|_\infty^{n-i} \|\mathbf{y}\|_\infty^i)$ combined with the fact that $\|\hat{\mathbf{U}} - \mathbf{U}\|_{\max} = o(1)$ to obtain

$$\begin{aligned} \max_{i,\ell} |\hat{\gamma}_i^\ell - \tilde{\gamma}_i^\ell| &\leq \max_{i,t,\ell} |\hat{U}_{i_1,t} \hat{U}_{i_2,t} \hat{U}_{i_3,t-|\ell|} \hat{U}_{i_4,t-|\ell|} - U_{i_1,t} U_{i_2,t} U_{i_3,t-|\ell|} U_{i_4,t-|\ell|}| \\ &= O(\|\hat{\mathbf{U}} - \mathbf{U}\|_{\max} \|\mathbf{U}\|_{\max}^3) \\ &= O_P[\eta[\psi^{-1}(nT)]^3] \end{aligned}$$

Therefore we conclude

$$\max_i |\hat{v}_i - \tilde{v}_i| \leq \max_{i,\ell} |\hat{\gamma}_i^\ell - \tilde{\gamma}_i^\ell| \sum_{|\ell| < T} |k(\ell/h)| = O_P\left(h\eta[\psi^{-1}(nT)]^3 \int |k(u)| du\right) = O_P(h\eta[\psi^{-1}(nT)]^3). \quad (\text{B.11})$$

The result then follows from the triangle inequality $\|\hat{\mathbf{Y}} - \mathbf{Y}\|_{\max} \leq \max_i |\hat{v}_i - \tilde{v}_i| + \max_i |\tilde{v}_i - v_i|$, expression (B.10) and (B.11). \square

Lemma B.14. *Let \mathbf{U}, \mathbf{V} be $T \times n$ matrices such that $\|\mathbf{U} - \mathbf{V}\|_{\max} \leq C_1$ and $\|\mathbf{V}\|_{\max} \leq C_2$, then*

$$\|\Sigma_{\mathbf{U}} - \Sigma_{\mathbf{V}}\|_{\max} \leq C_3 := C_1(2C_2 + C_1),$$

where $\Sigma_{\mathbf{U}} := \mathbf{U}'\mathbf{U}/T$ and $\Sigma_{\mathbf{V}} := \mathbf{V}'\mathbf{V}/T$. Furthermore, if $C_3 \leq \alpha\kappa(\Sigma_{\mathbf{V}}, \mathcal{S}, 3)/(|\mathcal{S}|(1 + \zeta)^2)$ for $\mathcal{S} \subseteq [n]$, $\zeta > 0$ and $\alpha \in [0, 1]$, then

$$(1 - \alpha)\kappa(\Sigma_{\mathbf{V}}, \mathcal{S}, \zeta) \leq \kappa(\Sigma_{\mathbf{U}}, \mathcal{S}, \zeta) \leq (1 + \alpha)\kappa(\Sigma_{\mathbf{V}}, \mathcal{S}, \zeta)$$

Proof. By the (reverse) triangle inequality we have $\|\mathbf{U}\|_{\max} - \|\mathbf{V}\|_{\max} \leq \|\mathbf{U} - \mathbf{V}\|_{\max}$, from which we conclude that $\|\mathbf{U}\|_{\max} \leq \|\mathbf{U} - \mathbf{V}\|_{\max} + \|\mathbf{V}\|_{\max} \leq C_1 + C_2$. Now $\|\Sigma_{\mathbf{U}} - \Sigma_{\mathbf{V}}\|_{\max} = \max_{1 \leq i, j \leq n} |T^{-1} \sum_{t=1}^T U_{it}U_{jt} - \mathbf{V}_{it}\mathbf{V}_{jt}|$

$V_{it}V_{ijt} \leq \max_{i,j,t} |U_{it}U_{jt} - V_{it}V_{jt}|$ and

$$\begin{aligned} |U_{it}U_{jt} - V_{it}V_{jt}| &\leq |(U_{it} - V_{it})U_{jt} + (U_{jt} - V_{jt})V_{it}| \leq \\ &\|\mathbf{U} - \mathbf{V}\|_{\max}(\|\mathbf{U}\|_{\max} + \|\mathbf{V}\|_{\max}) \leq C_1(2C_2 + C_1). \end{aligned}$$

For the second part of the lemma notice that for any $\mathbf{x} \in \mathbb{R}^n$ we have $|\mathbf{x}'\boldsymbol{\Sigma}_U\mathbf{x} - \mathbf{x}'\boldsymbol{\Sigma}_V\mathbf{x}| = |\mathbf{x}'(\boldsymbol{\Sigma}_U - \boldsymbol{\Sigma}_V)\mathbf{x}| \leq \|\boldsymbol{\Sigma}_U - \boldsymbol{\Sigma}_V\|_{\max}\|\mathbf{x}\|_1^2 \leq C_3\|\mathbf{x}\|_1^2$ by the first part. Also, if $\|\mathbf{x}_{S^c}\|_1 \leq \zeta\|\mathbf{x}_S\|_1$ we have that $\|\mathbf{x}\|_1 = \|\mathbf{x}_S\|_1 + \|\mathbf{x}_{S^c}\|_1 \leq (1 + \zeta)\|\mathbf{x}_S\|_1 \leq (1 + \zeta)\sqrt{\mathbf{x}'\boldsymbol{\Sigma}_V\mathbf{x}}|\mathcal{S}|/\kappa(\boldsymbol{\Sigma}_V, \mathcal{S}, \zeta)$ where the last inequality follows from the definition of compatibility condition. Thus $|\mathbf{x}'\boldsymbol{\Sigma}_U\mathbf{x} - \mathbf{x}'\boldsymbol{\Sigma}_V\mathbf{x}| \leq C_3(1 + \zeta)^2\mathbf{x}'\boldsymbol{\Sigma}_V\mathbf{x}|\mathcal{S}|/\kappa(\boldsymbol{\Sigma}_V, \mathcal{S}, \zeta) \leq \mathbf{x}'\boldsymbol{\Sigma}_V\mathbf{x}/2$, where the last inequality follows from the definition of compatibility condition. Therefore, we have that $(1 - \alpha)\mathbf{x}'\boldsymbol{\Sigma}_V\mathbf{x} \leq \mathbf{x}'\boldsymbol{\Sigma}_U\mathbf{x} \leq (1 + \alpha)\mathbf{x}'\boldsymbol{\Sigma}_V\mathbf{x}$ whenever $\|\mathbf{x}_{S^c}\|_1 \leq \zeta\|\mathbf{x}_S\|_1$. Take in infimum to conclude. \square

Lemma B.15. *Let $\mathbf{W} := (\mathbf{U}, \mathbf{V})$ and $\mathbf{Z} := (\mathbf{X}, \mathbf{Y})$ be $T \times (n+1)$ matrices such that $\|\mathbf{W} - \mathbf{Z}\|_{\max} \leq C_1$ and $\|\mathbf{Z}\|_{\max} \leq C_2$, then for any $\boldsymbol{\delta} \in \mathbb{R}^n$ we have*

$$\|\mathbf{U}'(\mathbf{V} - \mathbf{U}\boldsymbol{\delta})/T - \mathbf{X}'(\mathbf{Y} - \mathbf{X}\boldsymbol{\delta})/T\|_{\infty} \leq (1 + \|\boldsymbol{\delta}\|_1)C_1(2C_2 + C_1)$$

Proof. For convenience let $q := \mathbf{V} - \mathbf{U}\boldsymbol{\delta} \in \mathbb{R}^T$ and $r := \mathbf{Y} - \mathbf{X}\boldsymbol{\delta} \in \mathbb{R}^T$, then Hölder's inequality gives us $\|r\|_{\infty} \leq (1 + \|\boldsymbol{\delta}\|_1)\|Z\|_{\max} \leq (1 + \|\boldsymbol{\delta}\|_1)C_2$ and $\|q - r\|_{\infty} \leq (1 + \|\boldsymbol{\delta}\|_1)\|\mathbf{W} - \mathbf{Z}\|_{\max} \leq (1 + \|\boldsymbol{\delta}\|_1)C_1$. From the (reverse) triangle inequality we obtain $\|q\|_{\infty} \leq \|q - r\|_{\infty} + \|r\|_{\infty} \leq (1 + \|\boldsymbol{\delta}\|_1)(C_1 + C_2)$. Now, following the same steps in the proof of previous Lemma, we can upper bound the right hand side of the display by $\|\mathbf{U} - \mathbf{X}\|_{\max}\|q\|_{\infty} + \|q - r\|_{\infty}\|\mathbf{X}\|_{\max}$, which in turn can be upper bounded by the left hand side of the display. \square

Lemma B.16. *Under the same conditions of Theorems 1 and 2*

$$\begin{aligned} \|\nabla L(\boldsymbol{\theta}_1) - \nabla L_0(\boldsymbol{\theta}_1)\|_{\infty} &= O_P \left[\frac{\psi^{-1}(T)\psi^{-1}(nT)\psi^{-1}(n)\psi_{p/2}^{-1}(n)}{T^{1/4}} + \frac{\psi^{-1}(T)T^{1/4}}{\sqrt{n}} \right] \\ \|\nabla^2 L(\boldsymbol{\theta}) - \nabla^2 L_0(\boldsymbol{\theta})\|_{\max} &= O_P \left[\eta(n, T) [\psi^{-1}(nT) + \eta(n, T)] + \frac{\psi_{p/2}^{-1}(n^2)}{\sqrt{T}} \right], \end{aligned}$$

where $\nabla L_0(\boldsymbol{\theta}) := 2\mathbb{E}[\mathbf{U}_{-1t}(U_{1t} - \boldsymbol{\theta}'\mathbf{U}_{-1t})]$ and $\nabla^2 L_0(\boldsymbol{\theta}) := \mathbb{E}\mathbf{U}_{-1t}\mathbf{U}_{-1t}'$.

Proof. By the triangle inequality we have

$$\begin{aligned} \frac{1}{2} \|\nabla L(\boldsymbol{\theta}) - \nabla L_0(\boldsymbol{\theta})\|_\infty &= \|(\widehat{\mathbf{U}}_x - \mathbf{U}_x + \mathbf{U}_x)' \mathbf{V}/T - \mathbb{E}(\mathbf{U}'_x \mathbf{V}/T)\|_\infty \\ &\leq \|\mathbf{U}'_x \mathbf{V}/T - \mathbb{E}(\mathbf{U}'_x \mathbf{V}/T)\|_\infty + \|\widehat{\mathbf{U}}_x - \mathbf{U}_x\|_{\max} \|\mathbf{V}\|_\infty. \end{aligned}$$

Similarly, using Lemma 5.B

$$\begin{aligned} \|\nabla^2 L(\boldsymbol{\theta}) - \nabla^2 L_0(\boldsymbol{\theta})\|_{\max} &\leq \|\widehat{\mathbf{U}}'_x \widehat{\mathbf{U}}_x/T - \mathbf{U}'_x \mathbf{U}_x/T\|_{\max} + \|\mathbf{U}'_x \mathbf{U}_x/T - \mathbb{E}(\mathbf{U}'_x \mathbf{U}_x/T)\|_{\max} \\ &\leq \|\widehat{\mathbf{U}}_x - \mathbf{U}_x\|_{\max} (2\|\mathbf{U}_x\|_{\max} + \|\widehat{\mathbf{U}}_x - \mathbf{U}_x\|_{\max}) \\ &\quad + \|\mathbf{U}'_x \mathbf{U}_x/T - \mathbb{E}(\mathbf{U}'_x \mathbf{U}_x/T)\|_{\max}. \end{aligned}$$

By Corollary 1 and Assumption 3 we can bound in probability each of those terms

$$\begin{aligned} \|\mathbf{U}'_x \mathbf{V}/T - \mathbb{E}(\mathbf{U}'_x \mathbf{V}/T)\|_\infty &= O_P \left[\frac{\psi_{p/2}^{-1}(n)}{\sqrt{T}} \right] \\ \|\widehat{\mathbf{U}}_x - \mathbf{U}_x\|_{\max} &= O_P \left[\frac{\psi^{-1}(nT) \psi^{-1}(n) \psi_{p/2}^{-1}(n)}{T^{1/4}} + \frac{T^{1/4}}{\sqrt{n}} \right] =: O_P[\eta(n, T)] \\ \|\mathbf{V}\|_\infty &= \psi^{-1}(T) \\ \|\mathbf{U}_x\|_{\max} &= O_P[\psi^{-1}(nT)] \\ \|\mathbf{U}'_x \mathbf{U}_x/T - \mathbb{E}(\mathbf{U}'_x \mathbf{U}_x/T)\|_{\max} &= O_P \left[\frac{\psi_{p/2}^{-1}(n^2)}{\sqrt{T}} \right]. \end{aligned}$$

Therefore

$$\|\nabla L(\boldsymbol{\theta}) - \nabla L_0(\boldsymbol{\theta})\|_\infty = O_P \left[\frac{\psi_{p/2}^{-1}(n)}{\sqrt{T}} + \frac{\psi^{-1}(T) \psi^{-1}(nT) \psi^{-1}(n) \psi_{p/2}^{-1}(n)}{T^{1/4}} + \frac{\psi^{-1}(T) T^{1/4}}{\sqrt{n}} \right]$$

and

$$\|\nabla^2 L(\boldsymbol{\theta}) - \nabla^2 L_0(\boldsymbol{\theta})\|_{\max} = O_P \left[\eta(n, T) [\psi^{-1}(nT) + \eta(n, T)] + \frac{\psi_{p/2}^{-1}(n^2)}{\sqrt{T}} \right]$$

□

Lemma B.17. For $p > 0$, let $\psi_p : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ defined by $\psi_p(x) := x \mathbf{1}_{\{0 \leq x < t\}} + (\exp(x^p) - 1) \mathbf{1}_{\{x \geq t\}}$

where $t := \left(\frac{1-p}{p}\right)^{1/p}$. If $\|X\|_{\psi_p} < \infty$ then there exist constants $C_1 > 0$ and $C_2 > 0$ such that

$$\mathbb{P}(|X| > x) \leq C_1 \exp(-x^p/C_2) \quad x > 0.$$

In particular, if $0 < \|X\|_{\psi_p} < \infty$ we might take $C_1 = 2 + \mathbf{1}_{0 < p < 1} \exp((1-p)/p)$ and $C_2 = \|X\|_{\psi_p}^p$. Conversely, if there exist constants $C_1 > 0$ and $C_2 > 0$ such that $\mathbb{P}(|X| > x) \leq C_1 \exp(-x^p/C_2)$ for $x > 0$, then

$$\|X\|_{\psi_p} \leq \begin{cases} [(2C_1 + 1)C_2]^{1/p} \vee 2C_1 C_2^{1/p} p^{-1} \Gamma(1/p) & ; 0 < p < 1 \\ [(C_1 + 1)C_2]^{1/p} & ; p \geq 1, \end{cases}$$

where $\Gamma(\cdot)$ denotes the Gamma function.

Proof. If $\|X\|_{\psi_p} = 0$ then $X = 0$ a.s and the inequality holds for any choice of $C_1, C_2 > 0$. For the case when $0 < \|X\|_{\psi_p} < \infty$ we have by Markov inequality and the fact that $x \mapsto \exp a|x|^p$ is non-decreasing for $a > 0$

$$\mathbb{P}(|X| \geq x) = \mathbb{P}(\exp(a|X|^p) \geq \exp(ax^p)) \leq \exp(-ax^p) \mathbb{E} \exp(a|X|^p).$$

Also

$$\begin{aligned} \mathbb{E} \exp(a|X|^p) &= \mathbb{E} \exp(a|X|^p) \mathbf{1}_{a^{1/p}|X| < t} + \mathbb{E} \exp(a|X|^p) \mathbf{1}_{a^{1/p}|X| \geq t} \\ &\leq \exp((1-p)/p) \mathbf{1}_{0 < p < 1} + \mathbb{E} \psi_p(a^{1/p}|X|) + 1. \end{aligned}$$

Set $a = \|X\|_{\psi_p}^{-p}$ to conclude that the middle term is less or equal to 1.

For the converse we have for $a > 0$, by Fubini's Theorem

$$\begin{aligned} \mathbb{E} \exp(|aX|^p) - 1 &= \int \int_0^{|x|^p} a^p \exp(a^p y) dy \mathbb{P}(dx) \\ &= a^p \int_0^\infty \mathbb{P}(|X| \geq x^{1/p}) \exp(a^p x) dx. \end{aligned}$$

Since $\mathbb{P}(|X| > x) \leq C_1 \exp(-x^p/C_2)$, for $a < C_2^{-1/p}$, we have

$$\mathbb{E} \exp[(|aX|^p)] - 1 \leq a^p C_1 \int_0^\infty \exp\left[-x\left(\frac{1}{C_2} - a^p\right)\right] dx \leq \frac{a^p C_1}{C_2^{-1} - a^p}.$$

Also

$$\mathbb{E}|X| = \int_0^\infty \mathbb{P}(|X| > x) dx \leq C_1 \int_0^\infty \exp(-x^p/C_2) dx = \frac{C_1}{pC_2^{1/p}} \Gamma(1/p).$$

Therefore using the last two displays we have, for $0 < a < C_2^{-1/p}$

$$\begin{aligned} \mathbb{E}\psi_p(a|X|) &\leq \mathbb{E}|aX| \mathbf{1}_{0 < p < 1} + \mathbb{E} \exp[(|aX|)^p] - 1 \\ &\leq \frac{aC_1}{pC_2^{1/p}} \Gamma(1/p) \mathbf{1}_{0 < p < 1} + \frac{a^p C_1}{C_2^{1-a^p}}. \end{aligned}$$

When $p \geq 1$ the right hand side is less or equal than 1 for $a \leq [(1 + C_1)C_2]^{-1/p}$ hence $\|X\|_{\psi_p} \leq [(1 + C_1)C_2]^{1/p}$. For $0 < p < 1$, the right hand side is less or equal than 1 for $a \leq \{C_2^{-1/p}[(2C_1 + 1)^{-1/p} \wedge 2C_1 p^{-1} \Gamma(1/p)]\}^{-1}$ then $\|X\|_{\psi_p} \leq (2C_1 + 1)C_2^{1/p} \vee 2C_1 C_2^{1/p} p^{-1} \Gamma(1/p)$.

□

Lemma B.18. *For $p > 0$, there is a constant C_p only depending on p such that*

$$\|XY\|_{\psi_{p/2}} \leq C_p (\|X\|_{\psi_p} \vee \|Y\|_{\psi_p}),$$

where ψ_p defined as per Lemma (B.17).

Proof. If $\|X\|_{\psi_p} = 0$ or $\|Y\|_{\psi_p} = 0$ then $XY = 0$ a.s and the inequality hold trivially. Similarly if $\|X\|_{\psi_p} = \infty$ or $\|Y\|_{\psi_p} = \infty$. So we assume that $0 < \|X\|_{\psi_p} < \infty$ and $0 < \|Y\|_{\psi_p} < \infty$ and from Lemma (B.17) we have for $x > 0$

$$\begin{aligned} \mathbb{P}(|X| > x) &\leq K_p \exp[-(x/\|X\|_{\psi_p})^p] \\ \mathbb{P}(|Y| > x) &\leq K_p \exp[-(x/\|Y\|_{\psi_p})^p], \end{aligned}$$

where $K_p := 2 + \mathbf{1}_{0 < p < 1} \exp((1-p)/p)$. Then

$$\begin{aligned} \mathbb{P}(|XY| \geq x) &\leq \mathbb{P}(|X| \geq \sqrt{x}) + \mathbb{P}(|Y| \geq \sqrt{x}) \\ &\leq K_p \exp(-z^{p/2}/\|X\|_{\psi_p}^p) + K_p \exp(-z^{p/2}/\|Y\|_{\psi_p}^p) \\ &\leq 2K_p \exp(-z^{p/2}/D_p^p) \end{aligned}$$

where $D_p := \|X\|_{\psi_p} \vee \|Y\|_{\psi_p}$. Apply once again Lemma (B.17) to conclude that $\|XY\|_{\psi_{p/2}} \leq C_p D_p$

where $C_p = (2K_p + 1)^{1/p}$ for $p \geq 1$ otherwise $(4K_p + 1)^{1/p} \vee 8K_p \Gamma(2/p) p^{-1}$. □

Table 1: **Simulation Results: Size with $\phi = 0$.**

The table reports the empirical size of the test of remaining covariance structure. Panel (a) reports the case where the factors are known, whereas Panel (b) considers that the factors are unknown but the number of factors is known. Panels (c) and (d) present the results when the number of factors are determined, respectively, by the eigenvalue ratio test and the information criterion IC_1 . Factors are estimated by the usual principal component algorithm. Three nominal significance levels are considered: 0.01, 0.05, and 0.10. The table reports the results for the case where $\phi = 0$ in (5.3).

Panel(a): Known factors									
	<u>$T = 100$</u>			<u>$T = 500$</u>			<u>$T = 700$</u>		
	0.10	0.05	0.01	0.10	0.05	0.01	0.10	0.05	0.01
$n = 0.5 \times T$	0.08	0.03	0.01	0.10	0.05	0.01	0.09	0.04	0.01
$n = 1 \times T$	0.06	0.02	0.00	0.07	0.03	0.01	0.10	0.05	0.01
$n = 2 \times T$	0.07	0.02	0.00	0.07	0.02	0.00	0.08	0.04	0.00
$n = 3 \times T$	0.05	0.01	0.00	0.08	0.04	0.01	0.07	0.04	0.01

Panel(b): Known number of factors									
	<u>$T = 100$</u>			<u>$T = 500$</u>			<u>$T = 700$</u>		
	0.10	0.05	0.01	0.10	0.05	0.01	0.10	0.05	0.01
$n = 0.5 \times T$	0.23	0.13	0.02	0.14	0.06	0.02	0.11	0.05	0.01
$n = 1 \times T$	0.13	0.06	0.01	0.09	0.04	0.01	0.12	0.05	0.01
$n = 2 \times T$	0.09	0.04	0.01	0.08	0.04	0.01	0.09	0.04	0.00
$n = 3 \times T$	0.06	0.02	0.00	0.08	0.04	0.01	0.07	0.03	0.01

Panel(c): Information criterion (IC_1)									
	<u>$T = 100$</u>			<u>$T = 500$</u>			<u>$T = 700$</u>		
	0.10	0.05	0.01	0.10	0.05	0.01	0.10	0.05	0.01
$n = 0.5 \times T$	0.24	0.14	0.03	0.14	0.06	0.02	0.11	0.05	0.01
$n = 1 \times T$	0.14	0.07	0.02	0.09	0.04	0.01	0.12	0.05	0.01
$n = 2 \times T$	0.10	0.05	0.01	0.08	0.04	0.01	0.09	0.04	0.00
$n = 3 \times T$	0.07	0.03	0.01	0.08	0.04	0.01	0.07	0.03	0.01

Panel(d): Eigenvalue ratio									
	<u>$T = 100$</u>			<u>$T = 500$</u>			<u>$T = 700$</u>		
	0.10	0.05	0.01	0.10	0.05	0.01	0.10	0.05	0.01
$n = 0.5 \times T$	0.47	0.38	0.25	0.14	0.06	0.02	0.11	0.05	0.01
$n = 1 \times T$	0.14	0.07	0.02	0.09	0.04	0.01	0.12	0.05	0.01
$n = 2 \times T$	0.09	0.04	0.01	0.08	0.04	0.01	0.09	0.04	0.00
$n = 3 \times T$	0.06	0.02	0.00	0.08	0.04	0.01	0.07	0.03	0.01

Table 2: **Simulation Results: Size with $\phi = 0.5$.**

The table reports the empirical size of the test of remaining covariance structure. Panel (a) reports the case where the factors are known, whereas Panel (b) considers that the factors are unknown but the number of factors is known. Panels (c) and (d) present the results when the number of factors are determined, respectively, by the eigenvalue ratio test and the information criterion IC_1 . Factors are estimated by the usual principal component algorithm. Three nominal significance levels are considered: 0.01, 0.05, and 0.10. The table reports the results for the case where $\phi = 0.5$ in (5.3).

Panel(a): Known factors									
	<u>$T = 100$</u>			<u>$T = 500$</u>			<u>$T = 700$</u>		
	0.10	0.05	0.01	0.10	0.05	0.01	0.10	0.05	0.01
$n = 0.5 \times T$	0.09	0.03	0.01	0.11	0.06	0.01	0.10	0.05	0.01
$n = 1 \times T$	0.07	0.03	0.00	0.07	0.03	0.01	0.11	0.06	0.01
$n = 2 \times T$	0.08	0.02	0.00	0.08	0.03	0.00	0.09	0.05	0.00
$n = 3 \times T$	0.05	0.02	0.00	0.09	0.04	0.01	0.07	0.04	0.01

Panel(b): Known number of factors									
	<u>$T = 100$</u>			<u>$T = 500$</u>			<u>$T = 700$</u>		
	0.10	0.05	0.01	0.10	0.05	0.01	0.10	0.05	0.01
$n = 0.5 \times T$	0.24	0.14	0.02	0.15	0.07	0.02	0.12	0.06	0.01
$n = 1 \times T$	0.13	0.07	0.01	0.09	0.04	0.01	0.14	0.06	0.02
$n = 2 \times T$	0.09	0.04	0.01	0.08	0.04	0.01	0.09	0.05	0.00
$n = 3 \times T$	0.07	0.02	0.00	0.08	0.04	0.01	0.08	0.03	0.01

Panel(c): Information criterion (IC_1)									
	<u>$T = 100$</u>			<u>$T = 500$</u>			<u>$T = 700$</u>		
	0.10	0.05	0.01	0.10	0.05	0.01	0.10	0.05	0.01
$n = 0.5 \times T$	0.49	0.42	0.29	0.15	0.07	0.02	0.11	0.06	0.01
$n = 1 \times T$	0.15	0.09	0.02	0.10	0.04	0.01	0.14	0.06	0.01
$n = 2 \times T$	0.09	0.04	0.01	0.09	0.04	0.01	0.09	0.05	0.00
$n = 3 \times T$	0.07	0.03	0.00	0.10	0.04	0.01	0.08	0.03	0.01

Panel(d): Eigenvalue ratio									
	<u>$T = 100$</u>			<u>$T = 500$</u>			<u>$T = 700$</u>		
	0.10	0.05	0.01	0.10	0.05	0.01	0.10	0.05	0.01
$n = 0.5 \times T$	0.25	0.14	0.04	0.14	0.07	0.02	0.13	0.06	0.01
$n = 1 \times T$	0.15	0.07	0.02	0.10	0.04	0.01	0.13	0.06	0.02
$n = 2 \times T$	0.11	0.05	0.01	0.08	0.05	0.01	0.10	0.05	0.00
$n = 3 \times T$	0.08	0.03	0.01	0.09	0.04	0.01	0.08	0.03	0.01

Table 3: **Simulation Results: Power** ($\phi = 0$).

The table reports the empirical power of the test of remaining covariance structure. Panel (a) reports the case where the factors are known, whereas Panel (b) considers that the factors are unknown but the number of factors is known. Factors are estimated by the usual principal component algorithm. Three nominal significance levels are considered: 0.01, 0.05, and 0.10.

Panel(a): Known factors									
	<u>$T = 100$</u>			<u>$T = 500$</u>			<u>$T = 700$</u>		
	0.10	0.05	0.01	0.10	0.05	0.01	0.10	0.05	0.01
$n = 0.5 \times T$	1	1	1	1	1	1	1	1	1
$n = 1 \times T$	1	1	1	1	1	1	1	1	1
$n = 2 \times T$	1	1	1	1	1	1	1	1	1
$n = 3 \times T$	1	1	1	1	1	1	1	1	1

Panel(b): Known number of factors									
	<u>$T = 100$</u>			<u>$T = 500$</u>			<u>$T = 700$</u>		
	0.10	0.05	0.01	0.10	0.05	0.01	0.10	0.05	0.01
$n = 0.5 \times T$	0.33	0.18	0.03	0.99	0.97	0.83	0.99	0.99	0.94
$n = 1 \times T$	0.20	0.08	0.01	0.84	0.60	0.13	0.95	0.81	0.33
$n = 2 \times T$	0.16	0.07	0.01	0.83	0.55	0.11	0.94	0.82	0.34
$n = 3 \times T$	0.08	0.03	0.00	0.82	0.56	0.10	0.95	0.81	0.34

Panel(c): Eigenvalue ratio									
	<u>$T = 100$</u>			<u>$T = 500$</u>			<u>$T = 700$</u>		
	0.10	0.05	0.01	0.10	0.05	0.01	0.10	0.05	0.01
$n = 0.5 \times T$	0.15	0.09	0.02	0.99	0.97	0.83	0.99	0.99	0.94
$n = 1 \times T$	0.19	0.07	0.01	0.84	0.60	0.13	0.95	0.81	0.33
$n = 2 \times T$	0.16	0.07	0.01	0.83	0.55	0.11	0.94	0.82	0.34
$n = 3 \times T$	0.08	0.03	0.00	0.82	0.56	0.10	0.95	0.81	0.34

Panel(d): Information criterion (IC_1)									
	<u>$T = 100$</u>			<u>$T = 500$</u>			<u>$T = 700$</u>		
	0.10	0.05	0.01	0.10	0.05	0.01	0.10	0.05	0.01
$n = 0.5 \times T$	0.15	0.09	0.02	0.99	0.97	0.83	0.99	0.99	0.94
$n = 1 \times T$	0.19	0.07	0.01	0.84	0.60	0.13	0.95	0.81	0.33
$n = 2 \times T$	0.16	0.07	0.01	0.83	0.55	0.11	0.94	0.82	0.34
$n = 3 \times T$	0.08	0.03	0.00	0.82	0.56	0.10	0.95	0.81	0.34

Table 4: **Simulation Results: Power** ($\phi = 0.5$).

The table reports the empirical power of the test of remaining covariance structure. Panel (a) reports the case where the factors are known, whereas Panel (b) considers that the factors are unknown but the number of factors is known. Factors are estimated by the usual principal component algorithm. Three nominal significance levels are considered: 0.01, 0.05, and 0.10.

Panel(a): Known factors									
	<u>$T = 100$</u>			<u>$T = 500$</u>			<u>$T = 700$</u>		
	0.10	0.05	0.01	0.10	0.05	0.01	0.10	0.05	0.01
$n = 0.5 \times T$	1	1	1	1	1	1	1	1	1
$n = 1 \times T$	1	1	1	1	1	1	1	1	1
$n = 2 \times T$	1	1	1	1	1	1	1	1	1
$n = 3 \times T$	1	1	1	1	1	1	1	1	1

Panel(b): Known number of factors									
	<u>$T = 100$</u>			<u>$T = 500$</u>			<u>$T = 700$</u>		
	0.10	0.05	0.01	0.10	0.05	0.01	0.10	0.05	0.01
$n = 0.5 \times T$	0.36	0.18	0.03	1.00	1.00	0.91	1.00	1.00	1.00
$n = 1 \times T$	0.20	0.09	0.02	0.89	0.69	0.13	1.00	0.92	0.39
$n = 2 \times T$	0.18	0.07	0.01	0.98	0.59	0.13	1.00	0.96	0.36
$n = 3 \times T$	0.10	0.03	0.00	0.91	0.66	0.11	1.00	0.92	0.39

Panel(c): Eigenvalue ratio									
	<u>$T = 100$</u>			<u>$T = 500$</u>			<u>$T = 700$</u>		
	0.10	0.05	0.01	0.10	0.05	0.01	0.10	0.05	0.01
$n = 0.5 \times T$	0.15	0.11	0.03	1.00	1.00	0.96	1.00	1.00	1.00
$n = 1 \times T$	0.22	0.08	0.01	0.99	0.70	0.16	1.00	0.95	0.38
$n = 2 \times T$	0.17	0.07	0.01	0.94	0.62	0.12	0.98	0.88	0.37
$n = 3 \times T$	0.09	0.03	0.00	0.89	0.59	0.11	1.00	0.87	0.40

Panel(d): Information criterion (IC_1)									
	<u>$T = 100$</u>			<u>$T = 500$</u>			<u>$T = 700$</u>		
	0.10	0.05	0.01	0.10	0.05	0.01	0.10	0.05	0.01
$n = 0.5 \times T$	0.15	0.11	0.03	1.00	1.00	0.96	1.00	1.00	1.00
$n = 1 \times T$	0.22	0.08	0.01	0.99	0.70	0.16	1.00	0.95	0.38
$n = 2 \times T$	0.17	0.07	0.01	0.94	0.62	0.12	0.98	0.88	0.37
$n = 3 \times T$	0.09	0.03	0.00	0.89	0.59	0.11	1.00	0.87	0.40

Table 5: **Simulation Results: Informational Gains**

The table reports the average mean squared error (MSE) of three different prediction models over 5-fold cross-validation subsamples. The goal is to predict the first variable using information from the remaining $n - 1$. Panel (a) considers the case of Sparse Regression (SR) where Y_{1t} is LASSO-regressed on all the other variables. Panel (b) shows the results of Principal Component Regression (PCR). Finally, Panel (c) presents the results of **FarmPredict**. “N/A” means “not available”. Note that there is no factor selection for Sparse Regression. “Known Number” means that the number of factors is known.

Panel(a): Sparse Regression (SR)									
	<u>Known Number</u>			<u>Eigenvalue Ratio</u>			<u>Information Criterion (IC₁)</u>		
	$T = 100$	$T = 500$	$T = 700$	$T = 100$	$T = 500$	$T = 700$	$T = 100$	$T = 500$	$T = 700$
$n = 0.5 \times T$	0.57	0.35	0.34	N/A	N/A	N/A	N/A	N/A	N/A
$n = 1 \times T$	0.40	0.36	0.32	N/A	N/A	N/A	N/A	N/A	N/A
$n = 2 \times T$	0.39	0.33	0.31	N/A	N/A	N/A	N/A	N/A	N/A
$n = 3 \times T$	0.35	0.32	0.30	N/A	N/A	N/A	N/A	N/A	N/A

Panel(b): Principal Component Regression (PCR)									
	<u>Known Number</u>			<u>Eigenvalue Ratio</u>			<u>Information Criterion (IC₁)</u>		
	$T = 100$	$T = 500$	$T = 700$	$T = 100$	$T = 500$	$T = 700$	$T = 100$	$T = 500$	$T = 700$
$n = 0.5 \times T$	3.82	3.12	3.01	4.69	3.12	3.01	3.26	3.04	2.34
$n = 1 \times T$	3.09	2.35	2.34	4.05	3.35	3.34	3.22	3.02	2.32
$n = 2 \times T$	3.14	2.97	2.21	4.13	3.97	2.21	3.29	3.21	2.27
$n = 3 \times T$	3.83	3.00	2.33	3.83	3.00	2.33	3.12	3.00	2.28

Panel(c): FarmPredict									
	<u>Known Number</u>			<u>Eigenvalue Ratio</u>			<u>Information Criterion (IC₁)</u>		
	$T = 100$	$T = 500$	$T = 700$	$T = 100$	$T = 500$	$T = 700$	$T = 100$	$T = 500$	$T = 700$
$n = 0.5 \times T$	0.50	0.33	0.31	0.52	0.33	0.31	0.50	0.34	0.30
$n = 1 \times T$	0.32	0.29	0.28	0.37	0.29	0.28	0.53	0.28	0.27
$n = 2 \times T$	0.27	0.27	0.26	0.28	0.27	0.26	0.32	0.28	0.28
$n = 3 \times T$	0.22	0.21	0.21	0.22	0.21	0.21	0.34	0.27	0.27

Table 6: Forecasting Results.

The table reports the frequency each model is ranked the first, second, third and fourth best model among the four alternatives. Panel (a) considers the case when the factors are selected by the eigenvalue ratio procedure. Panel (b) presents the results when factors are selected by the information criterion IC_1 . Panels (c) and (d) consider the cases when the number of factors are pre-specified as either one or two. We present the results for each individual group of variables as well as for the full set of macroeconomic variables.

Panel (a): Optimal Factor Selection (eigenvalue ratio)																
Group	AR				SR				PCR				FarmPredict			
	1st	2nd	3rd	4th	1st	2nd	3rd	4th	1st	2nd	3rd	4th	1st	2nd	3rd	4th
(i) output and income	0.125	0.000	0.250	0.625	0.000	0.125	0.625	0.250	0.375	0.500	0.125	0.000	0.500	0.375	0.000	0.125
(ii) labor market	0.032	0.097	0.290	0.581	0.226	0.065	0.516	0.194	0.194	0.516	0.097	0.194	0.548	0.323	0.097	0.032
(iii) housing	0.100	0.100	0.300	0.500	0.400	0.400	0.100	0.100	0.000	0.200	0.400	0.400	0.500	0.300	0.200	0.000
(iv) consumption, orders and inventories	0.000	0.000	0.333	0.667	0.000	0.000	0.667	0.333	0.333	0.667	0.000	0.000	0.667	0.333	0.000	0.000
(v) money and credit	0.429	0.357	0.143	0.071	0.214	0.214	0.357	0.214	0.214	0.286	0.357	0.143	0.143	0.143	0.143	0.571
(vi) interest and exchange rates	0.368	0.211	0.263	0.158	0.526	0.316	0.158	0.000	0.053	0.263	0.211	0.474	0.053	0.211	0.368	0.368
(vii) prices	0.150	0.150	0.600	0.100	0.650	0.100	0.200	0.050	0.050	0.200	0.100	0.650	0.150	0.550	0.100	0.200
(viii) stock market	0.667	0.000	0.000	0.333	0.000	0.000	0.667	0.333	0.000	1.000	0.000	0.000	0.333	0.000	0.333	0.333
(ix) all	0.185	0.134	0.311	0.370	0.311	0.160	0.378	0.151	0.160	0.387	0.168	0.286	0.345	0.319	0.143	0.193
Panel (b): Optimal Factor Selection (IC_4)																
Group	AR				SR				PCR				FarmPredict			
	1st	2nd	3rd	4th	1st	2nd	3rd	4th	1st	2nd	3rd	4th	1st	2nd	3rd	4th
(i) output and income	0.125	0.125	0.188	0.563	0.063	0.250	0.500	0.188	0.250	0.375	0.313	0.063	0.563	0.250	0.000	0.188
(ii) labor market	0.032	0.097	0.258	0.613	0.226	0.032	0.548	0.194	0.226	0.581	0.065	0.129	0.516	0.290	0.129	0.065
(iii) housing	0.000	0.000	0.400	0.600	0.200	0.500	0.100	0.200	0.200	0.100	0.500	0.200	0.600	0.400	0.000	0.000
(iv) consumption, orders and inventories	0.000	0.000	0.333	0.667	0.167	0.000	0.500	0.333	0.167	0.667	0.167	0.000	0.667	0.333	0.000	0.000
(v) money and credit	0.571	0.286	0.071	0.071	0.143	0.429	0.357	0.071	0.143	0.286	0.429	0.143	0.143	0.000	0.143	0.714
(vi) interest and exchange rates	0.316	0.105	0.105	0.474	0.368	0.158	0.368	0.105	0.158	0.263	0.474	0.105	0.158	0.474	0.053	0.316
(vii) prices	0.100	0.150	0.650	0.100	0.500	0.300	0.150	0.050	0.100	0.150	0.100	0.650	0.300	0.400	0.100	0.200
(viii) stock market	0.667	0.000	0.000	0.333	0.000	0.667	0.333	0.000	0.000	0.333	0.667	0.000	0.333	0.000	0.000	0.667
(ix) all	0.176	0.118	0.277	0.429	0.252	0.227	0.378	0.143	0.176	0.353	0.269	0.202	0.395	0.303	0.076	0.227
Panel (c): Fixed Number of Factors ($r = 1$)																
Group	AR				SR				PCR				FarmPredict			
	1st	2nd	3rd	4th	1st	2nd	3rd	4th	1st	2nd	3rd	4th	1st	2nd	3rd	4th
(i) output and income	0.125	0.125	0.188	0.563	0.063	0.250	0.500	0.188	0.250	0.375	0.313	0.063	0.563	0.250	0.000	0.188
(ii) labor market	0.032	0.097	0.258	0.613	0.226	0.032	0.548	0.194	0.226	0.581	0.065	0.129	0.516	0.290	0.129	0.065
(iii) housing	0.000	0.000	0.400	0.600	0.200	0.500	0.100	0.200	0.200	0.100	0.500	0.200	0.600	0.400	0.000	0.000
(iv) consumption, orders and inventories	0.000	0.000	0.333	0.667	0.167	0.000	0.500	0.333	0.167	0.667	0.167	0.000	0.667	0.333	0.000	0.000
(v) money and credit	0.571	0.286	0.071	0.071	0.143	0.429	0.357	0.071	0.143	0.286	0.429	0.143	0.143	0.000	0.143	0.714
(vi) interest and exchange rates	0.316	0.105	0.105	0.474	0.368	0.158	0.368	0.105	0.158	0.263	0.474	0.105	0.158	0.474	0.053	0.316
(vii) prices	0.100	0.150	0.650	0.100	0.500	0.300	0.150	0.050	0.100	0.150	0.100	0.650	0.300	0.400	0.100	0.200
(viii) stock market	0.667	0.000	0.000	0.333	0.000	0.667	0.333	0.000	0.000	0.333	0.667	0.000	0.333	0.000	0.000	0.667
(ix) all	0.176	0.118	0.277	0.429	0.252	0.227	0.378	0.143	0.176	0.353	0.269	0.202	0.395	0.303	0.076	0.227
Panel (d): Fixed Number of Factors ($r = 2$)																
Group	AR				SR				PCR				FarmPredict			
	1st	2nd	3rd	4th	1st	2nd	3rd	4th	1st	2nd	3rd	4th	1st	2nd	3rd	4th
(i) output and income	0.063	0.125	0.250	0.563	0.063	0.063	0.625	0.250	0.250	0.625	0.063	0.063	0.625	0.188	0.063	0.125
(ii) labor market	0.065	0.129	0.226	0.581	0.226	0.097	0.516	0.161	0.226	0.452	0.097	0.226	0.484	0.323	0.161	0.032
(iii) housing	0.200	0.200	0.000	0.600	0.500	0.400	0.100	0.000	0.000	0.100	0.700	0.200	0.300	0.300	0.200	0.200
(iv) consumption, orders and inventories	0.167	0.167	0.167	0.500	0.167	0.167	0.500	0.167	0.333	0.333	0.333	0.000	0.333	0.333	0.000	0.333
(v) money and credit	0.500	0.357	0.071	0.071	0.214	0.357	0.143	0.286	0.143	0.286	0.429	0.143	0.143	0.000	0.357	0.500
(vi) interest and exchange rates	0.316	0.368	0.000	0.316	0.368	0.263	0.263	0.105	0.105	0.316	0.263	0.316	0.211	0.053	0.474	0.263
(vii) prices	0.100	0.100	0.100	0.700	0.500	0.150	0.250	0.100	0.200	0.150	0.550	0.100	0.200	0.600	0.100	0.100
(viii) stock market	0.667	0.000	0.000	0.333	0.000	0.667	0.333	0.000	0.000	0.333	0.667	0.000	0.333	0.000	0.000	0.667
(ix) all	0.193	0.193	0.126	0.487	0.286	0.202	0.361	0.151	0.176	0.345	0.311	0.168	0.345	0.261	0.202	0.193

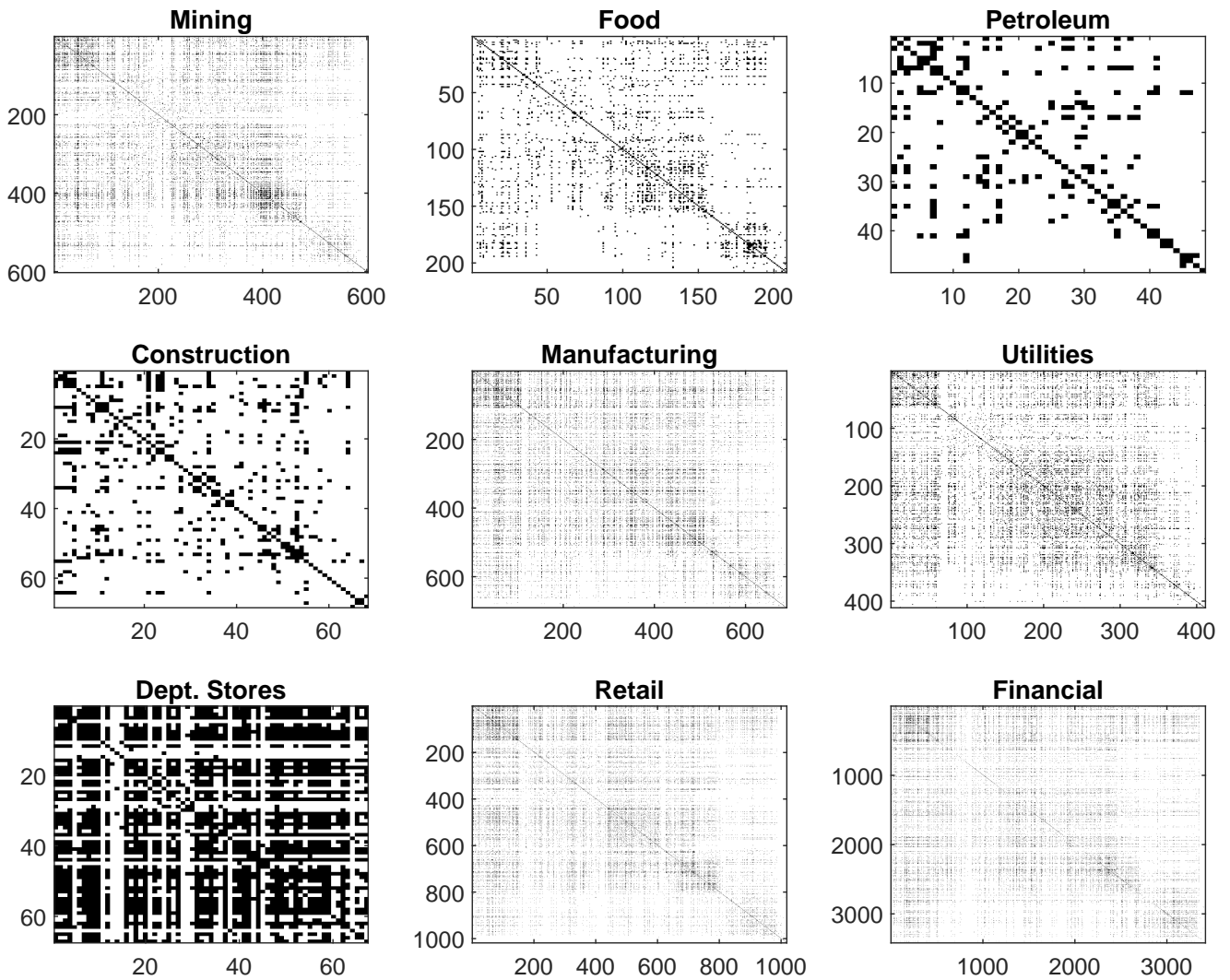


Figure 1: Correlations of returns larger than 0.15 in absolute value.

We estimate the correlations between all pairs of returns from a sample of nine specific sectors. The correlations that are higher than 0.15 in absolute value are shown as black dots in the figure. We consider the following sectors: mining, food, petroleum, construction, manufacturing, utilities, department stores, retail, and financial.

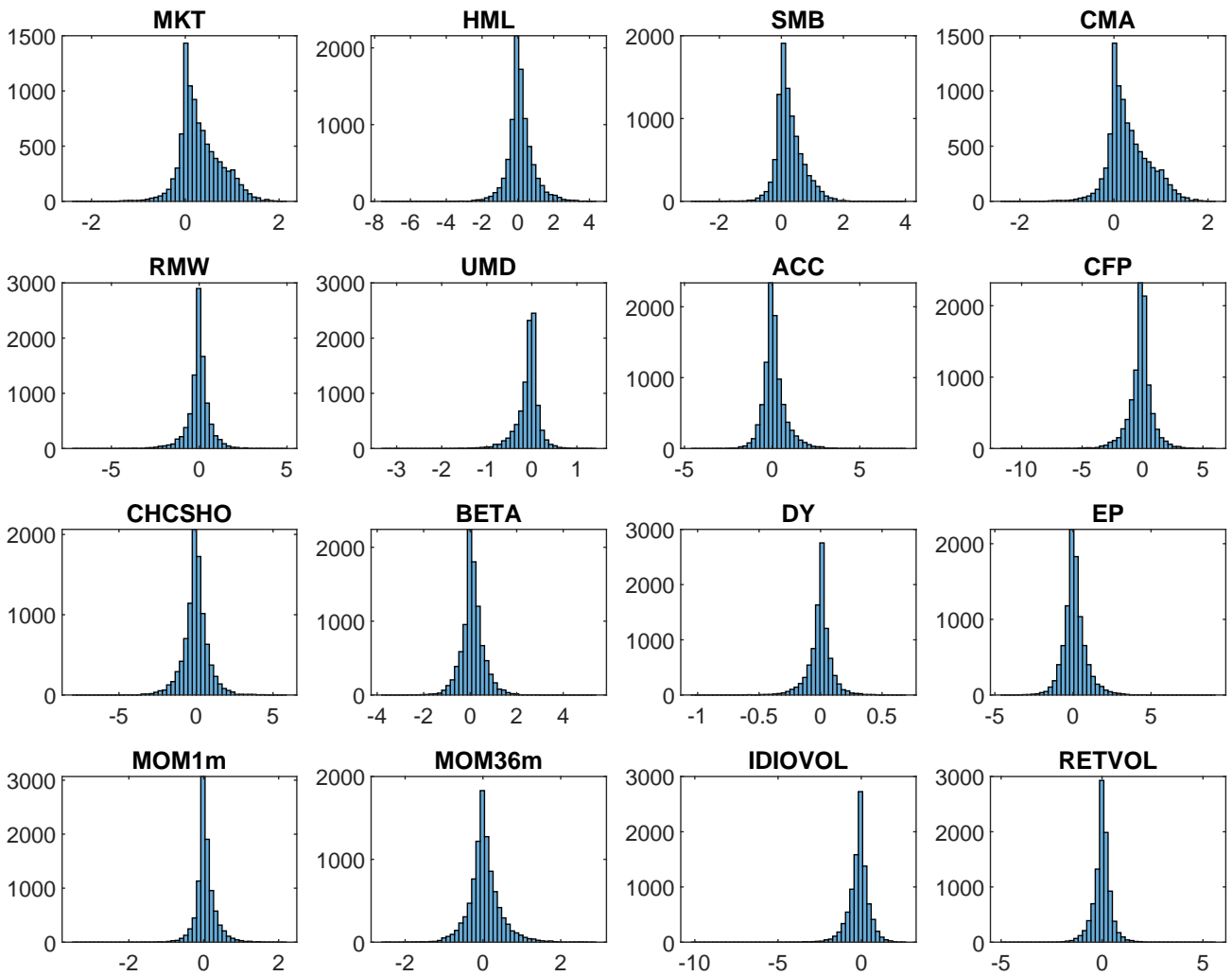


Figure 2: First-stage coefficient estimates.

The figure shows the empirical distribution of the first-stage regression where each excess returns are linearly regressed on 16 risk factors.

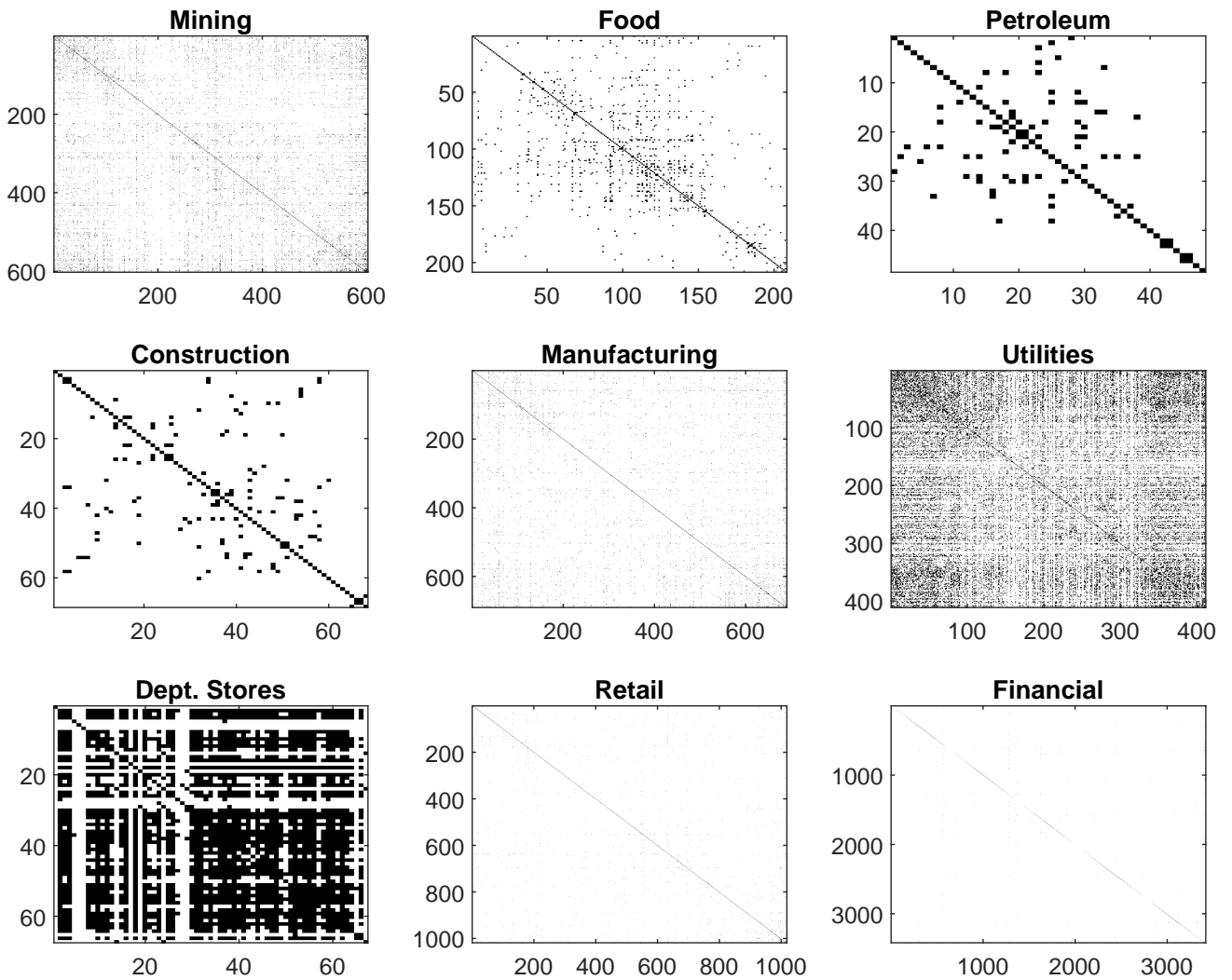


Figure 3: Correlations of first-stage residuals larger than 0.15 in absolute value. We estimate the correlations between all pairs of residuals from the first-stage OLS regression on 16 observed risk factors from a sample of nine specific sectors. The correlations that are higher than 0.15 in absolute value are shown as black dots in the figure. We consider the following sectors: mining, food, petroleum, construction, manufacturing, utilities, department stores, retail, and financial.

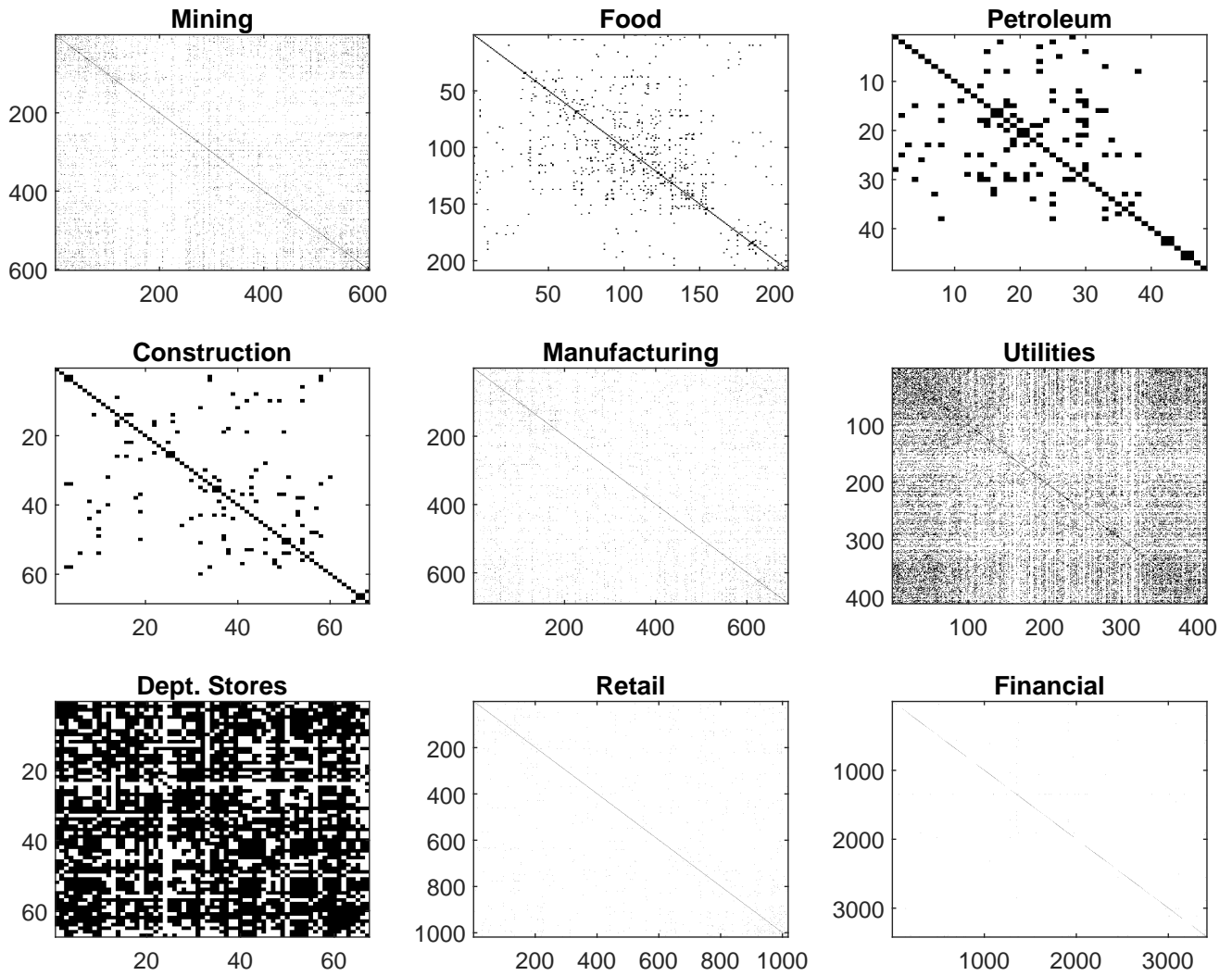


Figure 4: Correlations of second-stage residuals larger than 0.15 in absolute value. We estimate the correlations between all pairs of residuals from the second-stage principal component analysis from a sample of nine specific sectors. The correlations that are higher than 0.15 in absolute value are shown as black dots in the figure. We consider the following sectors: mining, food, petroleum, construction, manufacturing, utilities, department stores, retail, and financial.

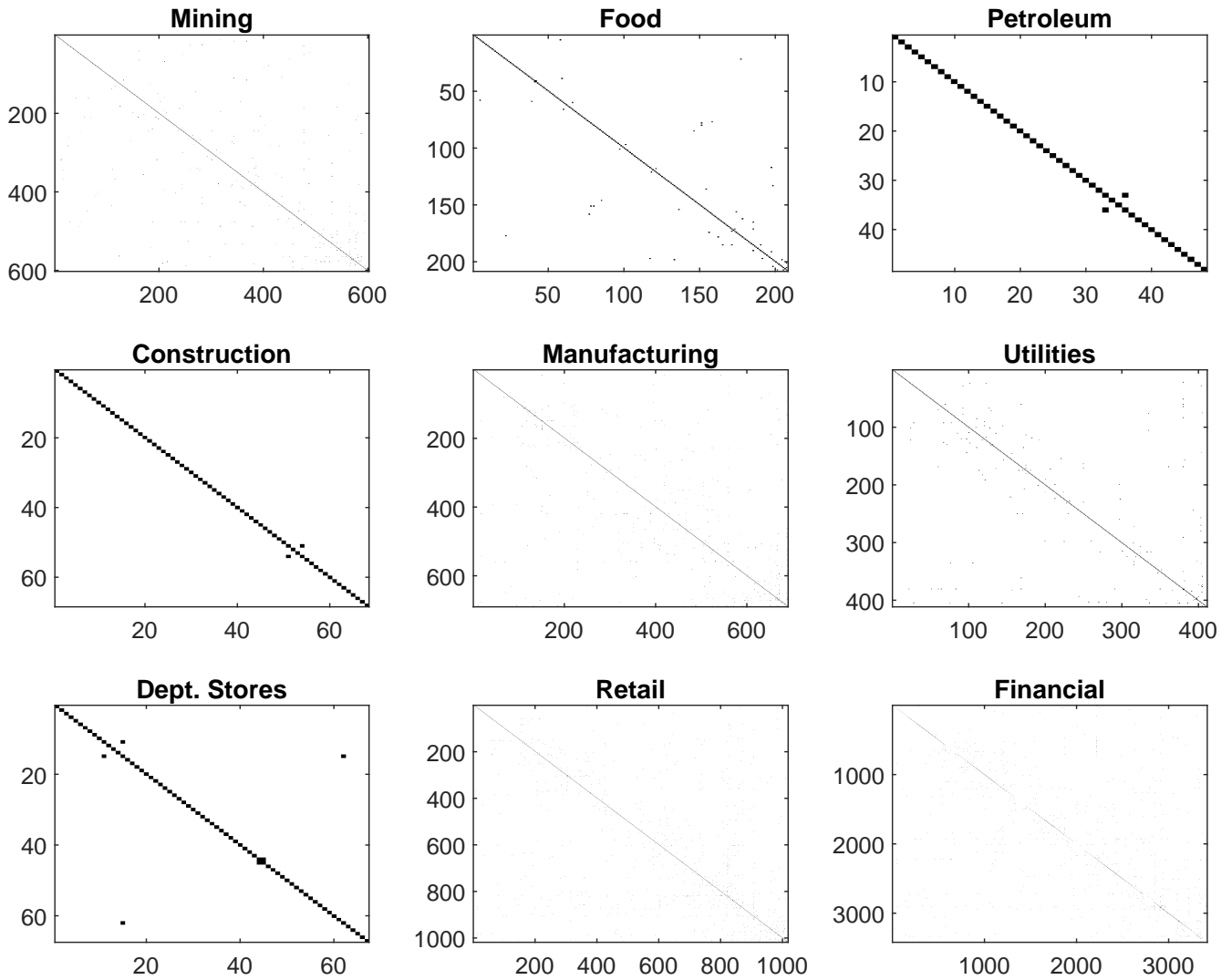


Figure 5: Partial correlations of second-stage residuals larger than 0.15 in absolute value. We estimate the partial correlations between all pairs of residuals from the second-stage LASSO regression from a sample of nine specific sectors. The correlations that are higher than 0.15 in absolute value are shown as black dots in the figure. We consider the following sectors: mining, food, petroleum, construction, manufacturing, utilities, department stores, retail, and financial.

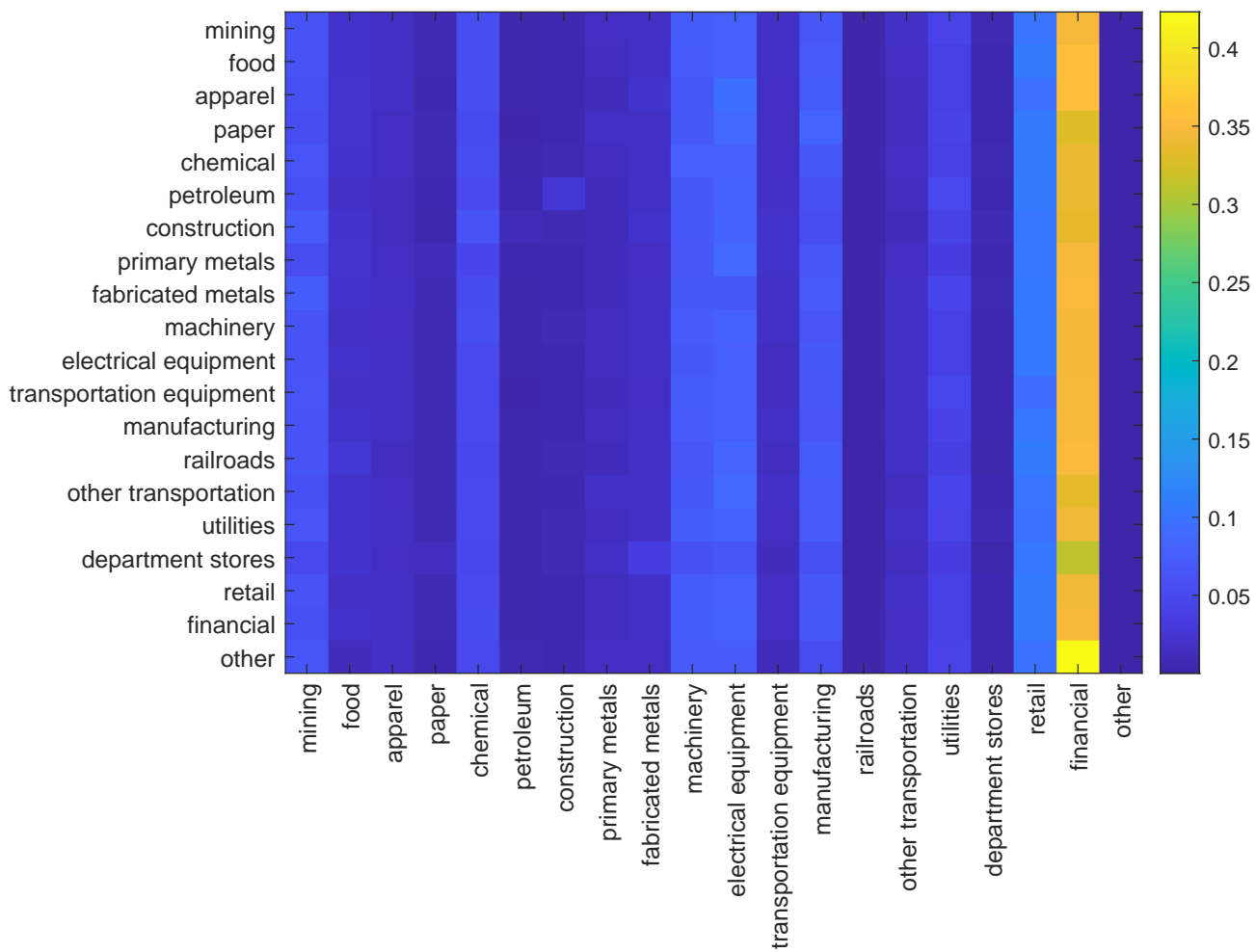


Figure 6: Variable Selection Frequency.

We report how often that variables from column sectors are selected in the third-stage LASSO regression for firms on line sectors . The numbers are normalized by the total number of firms in each sector.

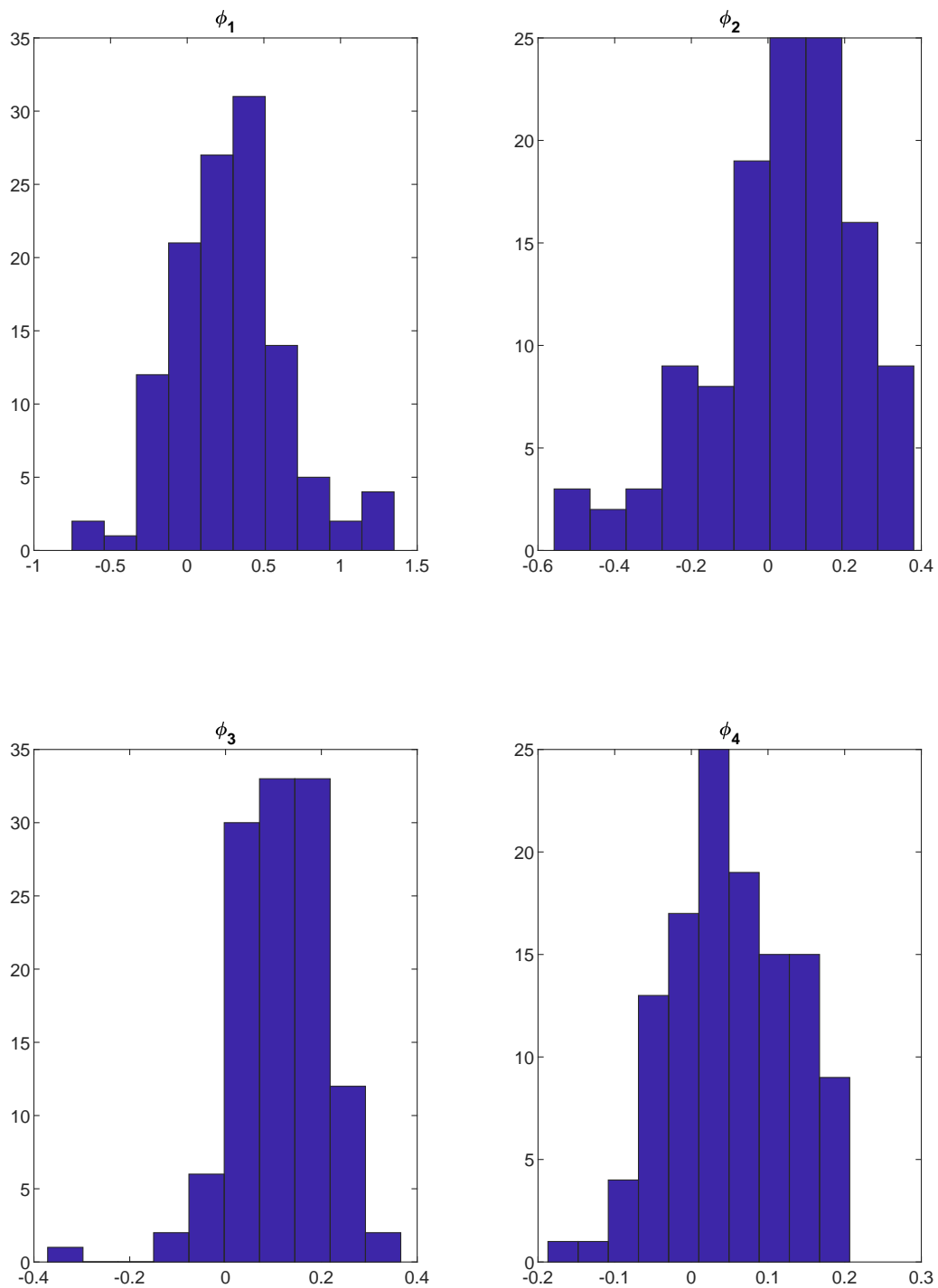


Figure 7: AR coefficient estimates.

The figure illustrates the empirical distribution of the ordinary least squares (OLS) estimation of the coefficients of an fourth-order autoregressive, AR(4), model across the 119 macroeconomic time series. Each panel relates to one specific coefficient.

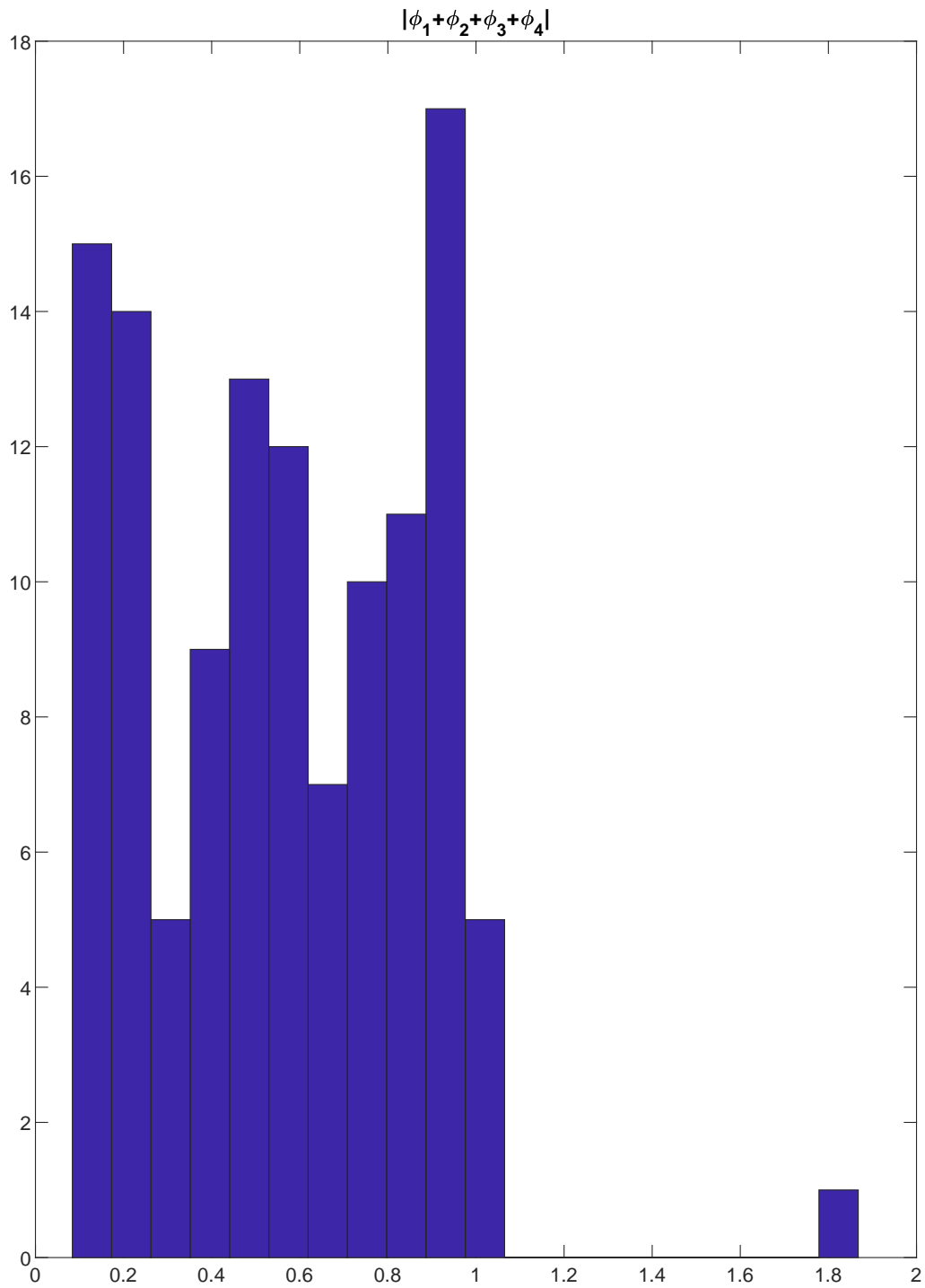


Figure 8: Absolute sum of AR coefficient estimates.

The figure illustrates the empirical distribution of the absolute sum of the ordinary least squares (OLS) estimation of the coefficients of an fourth-order autoregressive, AR(4), model across the 119 macroeconomic time series.

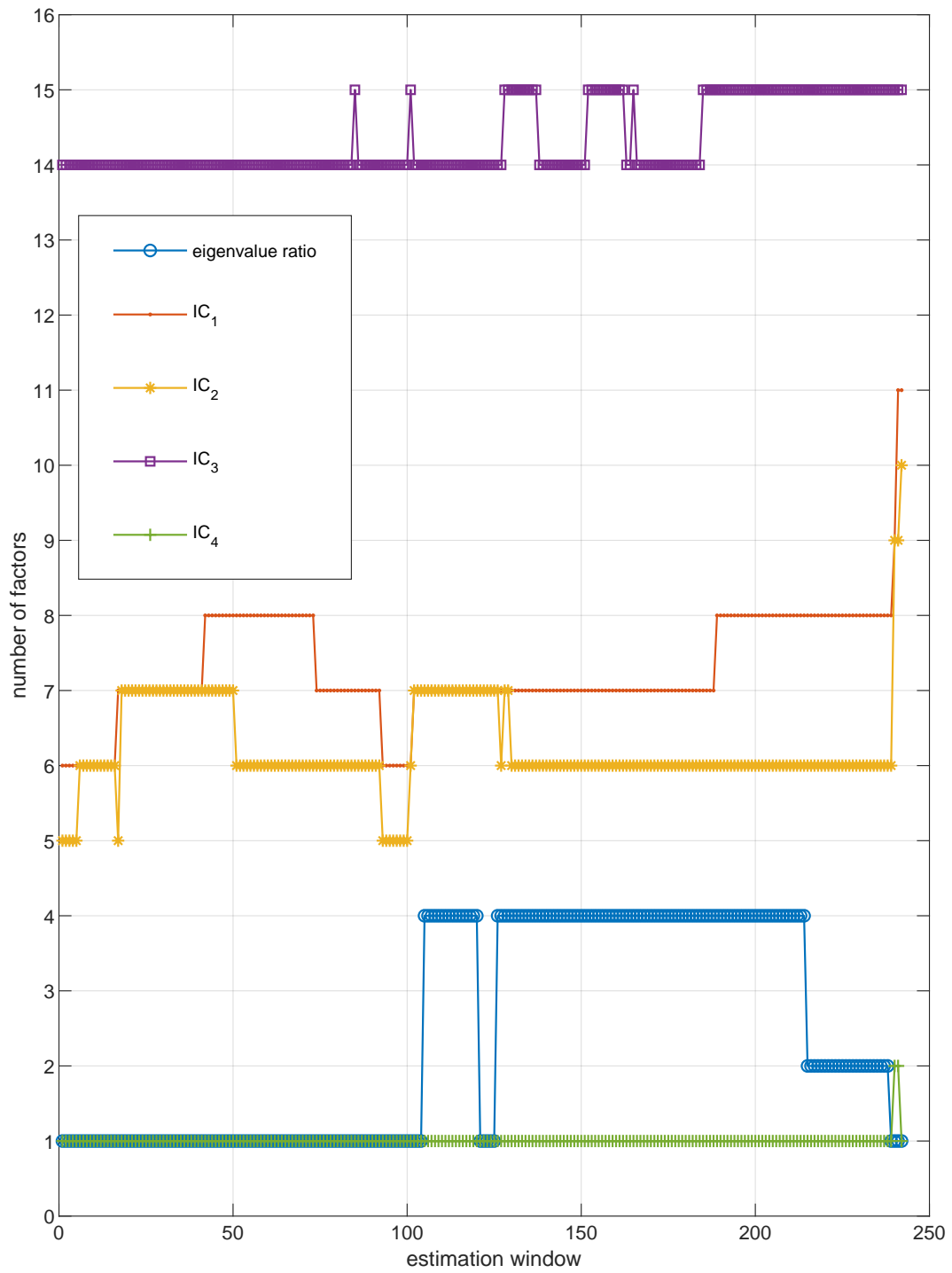


Figure 9: Estimated number of factors.

The figure illustrates the number of selected factors over the estimation windows. The figure reports the results for the eigenvalue ratio procedure and the four information criteria discussed in the paper.

References

- Abadie, A., A. Diamond, and J. Hainmueller (2010). Synthetic control methods for comparative case studies: Estimating the effect of California’s tobacco control program. *Journal of the American Statistical Association* 105, 493–505.
- Abadie, A. and J. Gardeazabal (2003). The economic costs of conflict: A case study of the Basque country. *American Economic Review* 93, 113–132.
- Andreou, E. and E. Ghysels (2021). Predicting the VIX and the volatility risk premium: The role of short-run funding spreads volatility factors. *Journal of Econometrics* 220, 366–398.
- Andrews, D. W. K. (1991). Heteroskedasticity and autocorrelation consistent covariance matrix estimation. *Econometrica* 59(3), 817–858.
- Bai, J. (2003). Inferential theory for factor models of large dimensions. *Econometrica* 71, 2135–171.
- Bai, J. and Y. Liao (2017). Inferences in panel data with interactive effects using large covariance matrices. *Journal of Econometrics* 200, 59–78.
- Bai, J. and S. Ng (2002). Determining the number of factors in approximate factor models. *Econometrica* 70, 191–221.
- Bai, J. and S. Ng (2003). Inferential theory for factor models of large dimensions. *Econometrica* 71, 135–171.
- Bai, J. and S. Ng (2006). Confidence intervals for diffusion index forecasts and inference for factor augmented regressions. *Econometrica* 74, 1133–1155.
- Barigozzi, M. and C. Brownlees (2019). NETS: Network estimation for time series. *Journal of Applied Econometrics* 34, 347–364.
- Barigozzi, M. and M. Hallin (2016). Generalized dynamic factor models and volatilities: Recovering the market volatility shocks. *Econometrics Journal* 19, C33–C60.
- Barigozzi, M. and M. Hallin (2017a). Generalized dynamic factor models and volatilities: Estimation and forecasting. *Journal of Econometrics* 201, 307–321.
- Barigozzi, M. and M. Hallin (2017b). A network analysis of the volatility of high-dimensional financial series. *Journal of the Royal Statistical Society - series C* 66, 581–605.

- Barigozzi, M. and M. Hallin (2020). Generalized dynamic factor models and volatilities: consistency, rates, and prediction intervals. *Journal of Econometrics* 116, 4–34.
- Barigozzi, M., M. Hallin, S. Soccorsi, and R. von Sachs (2020). Time-varying general dynamic factor models and the measurement of financial connectedness. *Journal of Econometrics*. forthcoming.
- Bernanke, B., J. Boivin, and P. Elias (2005). Measuring the effects of monetary policy: A factor-augmented vector autoregressive (FAVAR) approach. *The Quarterly Journal of Economics* 120, 387–422.
- Brito, D., M. Medeiros, and R. Ribeiro (2018). Forecasting large realized covariance matrices: The benefits of factor models and shrinkage. Technical Report 3163668, SSRN.
- Brownlees, C., G. Gudmundsson, and G. Lugosi (2020). Community detection in partial correlation network models. *Journal of Business & Economic Statistics*. forthcoming.
- Cai, T. (2017). Global testing and large-scale multiple testing for high-dimensional covariance structures. *Annual Review of Statistics and its Application* 4, 4.1–4.24.
- Cai, T. and Z. Ma (2013). Optimal hypothesis testing for high dimensional covariance matrices. *Bernoulli* 19, 2359–2388.
- Cai, T., Z. Ren, and H. Zhou (2016). Estimating structured high-dimensional covariance and precision matrices: Optimal rates and adaptive estimation. *Electronic Journal of Statistics* 10, 1–59.
- Carvalho, C., R. Masini, and M. Medeiros (2018). Arco: An artificial counterfactual approach for high-dimensional panel time-series data. *Journal of Econometrics* 207, 352–380.
- Chen, S., L.-X. Zhang, and P.-S. Zhong (2010). Tests for high-dimensional covariance matrices. *Journal of the American Statistical Association* 105, 810–819.
- Chernozhukov, V., D. Chetverikov, and K. Kato (2013). Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors. *Annals of Statistics* 41, 2786–2819.
- Chernozhukov, V., D. Chetverikov, and K. Kato (2018). Inference on causal and structural parameters using many moment inequalities.

- Diebold, F. and K. Yilmaz (2014). On the network topology of variance decompositions: Measuring the connectedness of financial firms. *Journal of Econometrics* 182, 119–134.
- Doukhan, P. and S. Louhichi (1999). A new weak dependence condition and applications to moment inequalities. *Stochastic Processes and their Applications* 84, 313–342.
- Fama, E. and K. French (1993). Common risk factors in the returns on stocks and bonds. *Journal of Financial Economics* 33, 3–56.
- Fama, E. and K. French (2015). A five-factor asset pricing model. *Journal of Financial Economics* 116, 1–22.
- Fan, J., Y. Fan, and J. Lv (2008). High dimensional covariance matrix estimation using a factor model. *Journal of Econometrics* 147, 186–197.
- Fan, J., Y. Ke, and K. Wang (2020). Factor-adjusted regularized model selection. *Journal of Econometrics* 216, 71–85.
- Fan, J., Q. Li, and Y. Wang (2017). Estimation of high dimensional mean regression in the absence of symmetry and light tail assumptions. *Journal of the Royal Statistical Society: Series B* 79, 247–265.
- Fan, J., R. Li, C.-H. Zhang, and H. Zou (2020). *Statistical Foundations of Data Science*. CRC Press.
- Fan, J., Y. Liao, and M. Mincheva (2013). Large covariance estimation by thresholding principal orthogonal complements. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 75, 603–680.
- Fan, J., R. Masini, and M. Medeiros (2020). Do we exploit all information for counterfactual analysis? benefits of factor models and idiosyncratic correction. Working paper, Princeton University.
- Feng, G., S. Giglio, and D. Xiu (2020). Taming the factor zoo: A test of new factors. *Journal of Finance* 75, 1327–1370.
- Gagliardini, P., E. Ossola, and P. Scaillet (2019). A diagnostic criterion for approximate factor structure. *Journal of Econometrics* 212, 503–521.

- Gagliardini, P., E. Ossola, and P. Scaillet (2020). Estimation of large dimensional conditional factor models in finance. In S. Durlauf, L. Hansen, J. Heckman, and R. Matzkin (Eds.), *Handbook of Econometrics*, Volume Volume 7A, pp. 219–282.
- Giannone, D., M. Lenza, and G. Primiceri (2018). Economic predictions with big data: The illusion of sparsity. Working paper, Northwestern University.
- Giessing, A. and J. Fan (2020). Bootstrapping ℓ_p -statistics in high dimensions.
- Giglio, S. and D. Xiu (2020). Asset pricing with omitted factors. *Journal of Political Economy*. forthcoming.
- Gobillon, L. and T. Magnac (2016). Regional policy evaluation: Interactive fixed effects and synthetic controls. *Review of Economics and Statistics* 98, 535–551.
- Gu, S., B. Kelly, and D. Xiu (2020). Empirical asset pricing via machine learning. *Review of Financial Studies* 33, 2223–2273.
- Guo, X. and C. Tang (2020). Specification tests for covariance structures in high-dimensional statistical models. *Biometrika*. forthcoming.
- Horenstein, S. A. A. (2013). Eigenvalue ratio test for the number of factors. *Econometrica* 81, 1203–1227.
- Kock, A. and L. Callot (2015). Oracle inequalities for high dimensional vector autoregressions. *Journal of Econometrics* 186, 325–344.
- Lam, C. and J. Fan (2009). Sparsistency and rates of convergence in large covariance matrix estimation. *Annals of Statistics* 37, 4254–4278.
- Ledoit, O. and M. Wolf (2002). Some hypothesis tests for the covariance matrix when the dimension is large compared to the sample size. *Annals of Statistics* 30, 1081–1102.
- Ledoit, O. and M. Wolf (2004). A well-conditioned estimator for large-dimensional covariance matrices. *Journal of Multivariate Analysis* 88, 365–411.
- Ledoit, O. and M. Wolf (2012). Nonlinear shrinkage estimation of large-dimensional covariance matrices. *Annals of Statistics* 40, 1024–1060.

- Ledoit, O. and M. Wolf (2017). Nonlinear shrinkage of the covariance matrix for portfolio selection: Markowitz meets Goldilocks. *Review of Financial Studies* 30, 4349–4388.
- Ledoit, O. and M. Wolf (2020). Analytical nonlinear shrinkage of large-dimensional covariance matrices. *Annals of Statistics*. forthcoming.
- Ledoit, O. and M. Wolf (2021a). The power of (non-)linear shrinking: A review and guide to covariance matrix estimation. *Journal of Financial Econometrics*. forthcoming.
- Ledoit, O. and M. Wolf (2021b). Quadratic shrinkage for large covariance matrices. *Bernoulli*. forthcoming.
- Li, W. and Y. Qin (2014). Hypothesis testing for high-dimensional covariance matrices. *Journal of Multivariate Analysis* 128, 108–119.
- Masini, R. and M. Medeiros (2019). Counterfactual analysis with artificial controls: Inference, high dimensions and nonstationarity. Working Paper 3303308, SSRN.
- Masini, R., M. Medeiros, and E. Mendes (2019). Regularized estimation of high-dimensional vector autoregressions with weakly dependent innovations. Technical Report 1912.09002, arxiv.
- McCracken, M. and S. Ng (2016). FRED-MD: A monthly database for macroeconomic research. *Journal of Business & Economic Statistics* 34, 574–589.
- Medeiros, M. and E. Mendes (2016). ℓ_1 -regularization of high-dimensional time-series models with non-gaussian and heteroskedastic errors. *Journal of Econometrics* 191, 255–271.
- Merlevède, F., M. Peligrad, and E. Rio (2009). Bernstein inequality and moderate deviations under strong mixing conditions. In C. Houdré, V. Koltchinskii, D. Mason, and M. Peligrad (Eds.), *High Dimensional Probability V: The Luminy Volume*, Volume Volume 5, pp. 273–292. Institute of Mathematical Statistics.
- Moon, R. and M. Weidner (2015). Linear regression for panel with unknown number of factors as interactive fixed effects. *Econometrica* 83, 1543–1579.
- Moskowitz, T. and M. Grinblatt (1999). Do industries explain momentum? *Journal of Finance* 54, 1249–1290.

- Negahban, S., P. Ravikumar, M. Wainwright, and B. Yu (2012). A unified framework for high-dimensional analysis of m -estimators with decomposable regularizers. *Statistical Science* 27, 538–557.
- Onatski, A., M. Moreira, and M. Hallin (2013). Asymptotic power of sphericity tests for high-dimensional data. *Annals of Statistics* 41, 1204–1231.
- Rio, E. (1994). Inégalités de moments pour les suites stationnaires et fortement mélangées. *Comptes rendus Acad. Sci. Paris, Série I* 318, 355–360.
- Stock, J. and M. Watson (2002a). Forecasting using principal components from a large number of predictors. *Journal of the American Statistical Association* 97, 1167–1179.
- Stock, J. and M. Watson (2002b). Macroeconomic forecasting using diffusion indexes. *Journal of Business & Economic Statistics* 20, 147–162.
- van de Geer, S. and P. Bühlmann (2009). On the conditions used to prove oracle results for the lasso. *Electronic Journal of Statistics* 3, 1360–1392.
- van der Vaart, A. and J. Wellner (1996). *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer.
- Zheng, S., Z. Chen, H. Cui, and R. Li (2019). Hypothesis testing on linear structures of high-dimensional covariance matrix. *Annals of Statistics* 47, 3300–3334.
- Zheng, S., G. Cheng, J. Guo, and H. Zhu (2019). Test for high-dimensional correlation matrices. *Annals of Statistics* 47, 2887–2921.