# High-Dimensional Spatial Arbitrage Pricing Theory with Heterogeneous Interactions

Zhaoxing Gao[1], Sihan Tu[2], and Ruey Tsay[3*]

[1]School of Mathematical Sciences, University of Electronic Science and Technology of China

[2]School of Management, Zhejiang University

[3]Booth School of Business, University of Chicago

## Abstract

This paper investigates estimation and inference of a Spatial Arbitrage Pricing Theory (SAPT) model that integrates spatial interactions with multi-factor analysis, accommodating both observable and latent factors. Building on the classical mean-variance analysis, we introduce a class of Spatial Capital Asset Pricing Models (SCAPM) that account for spatial effects in high-dimensional assets, where we define *spatial rho* as a counterpart to market beta in CAPM. We then extend SCAPM to a general SAPT framework under a *complete* market setting by incorporating multiple factors. For SAPT with observable factors, we propose a generalized shrinkage Yule-Walker (SYW) estimation method that integrates ridge regression to estimate spatial and factor coefficients. When factors are latent, we first apply an autocovariance-based eigenanalysis to extract factors, then employ the SYW method using the estimated factors. We establish asymptotic properties for these estimators under high-dimensional settings where both the dimension and sample size diverge. Finally, we use simulated and real data examples to demonstrate the efficacy and usefulness of the proposed model and method.

**Keywords:** Spatial Arbitrage Pricing Theory, Multi-factor Analysis, Yule-Walker Estimation, Eigenanalysis, High Dimension

---

*Corresponding author: ruey.tsay@chicagobooth.edu (R.S. Tsay). Booth School of Business, University of Chicago, 5807 South Woodlawn Avenue, Chicago, IL, 60637, USA.

# 1 Introduction

With the rapid advancement in information technology, large-scale datasets have become ubiquitous across all scientific areas with important applications. These datasets also introduce new analytical challenges in financial econometrics and statistics, particularly in high-dimensional settings. As a fundamental tool for dimension reduction and feature extraction, factor models provide a crucial link between economic theory and data analysis. Since the seminal work of Markowitz (1952) on portfolio theory, factor-based pricing models have played a central role in asset pricing, investment analysis and risk assessment. The Capital Asset Pricing Model (CAPM), developed by Sharpe, Lintner, and Mossin in the 1960s, introduced the concept of *market beta* to quantify systematic risk-return relationships. Ross (1976) proposed the Arbitrage Pricing Theory (APT), which extended the single-factor CAPM by incorporating multiple systematic risk factors under no-arbitrage principles, allowing for a more flexible representation of expected returns. Modern factor-based pricing research has evolved into two dominant approaches to address the growing market complexity. The first approach, developed by Fama and French (1993, 2015), relies on the theory-driven observable factors, such as market returns and firm characteristics. Building on this framework, numerous factor models for asset returns have been proposed; for instance, Feng et al. (2020) propose the Double-Selection LASSO to evaluate the marginal contribution of individual factors relative to an existing high-dimensional factor set. While these models offer strong economic interpretability, their fixed factor structures limit their ability to capture modern dynamic market interactions. Recent studies by Forni et al. (2000), Bai and Ng (2002), Bai (2003), Forni et al. (2005), Lam and Yao (2012), Fan et al. (2013), Gao and Tsay (2022, 2023), among others, have focused on latent factor models as an alternative approach. These models provide a methodology for inferring unobserved common factors from covariance structures. Lettau and Pelger (2020) and Giglio et al. (2025) further demonstrate the effectiveness of their tailored latent factor models in asset pricing, offering deeper insights into the underlying structure of financial markets. Liu et al. (2025) show that one can improve the estimation of portfolio risk by augmenting the Fama and French factors with latent factors extracted from a matrix-variate dataset of asset returns.

Despite the effectiveness of factor models in explaining cross-sectional and dynamic dependence, many economic and financial applications often manifest intricate spatial interconnections. Consider, for example, the spatial distribution of economic indicators across regions, where the performance of one region may influence its neighbors; see Anselin (1988) and Cressie (2015). Since the seminal work of Cliff and Ord (1973) on spatial autocorrelations, spatial models are often used to model cross-sectional dependence of different economic units or individuals at different locations. More recently, the spatial models have been extended to spatial dynamic panel data (SDPD) models by adding a time-lagged direction to account for serial correlations across different economic units or individuals; see, for example, Lee and Yu (2010). Empirically, the spatial interactions among the panel may exist in many large-dimensional economic and financial systems, together with other comovements or common factors. For example, Pirinsky and Wang (2006) found the spatial effect in the U.S. equity market by studying the comovements of common stock returns of U.S. corpora-

tions in the same geographic area; Kou et al. (2018) proposed an asset pricing model with spatial interactions and discovered significant spatial interactions in the futures contracts on S&P/Case-Shiller Home Price Indices. Therefore, augmenting factor models with spatial interactions not only extends these models with additional common factors but also enriches spatio-temporal models by integrating common factor structures.

In this paper, we focus on spatial panel models with common factors in the context of arbitrage pricing under high-dimensional settings. Building on the classical mean-variance analysis, we first introduce a class of Spatial Capital Asset Pricing Models (SCAPM) that account for spatial effects in high-dimensional assets under a "*complete market*" or "*minimum complete market*" assumption, where we introduce a *spatial rho* as a counterpart to market beta in CAPM. Within the spatial CAPM framework, we extend the model to a Spatial Arbitrage Pricing Theory (SAPT) by incorporating a multifactor structure. This formulation captures both systematic risk factors and spatial spillover effects, offering a unified approach to modeling interdependencies in asset returns.

While prior studies, such as Pesaran and Tosetti (2011), Kou et al. (2018), Bai and Li (2021), Yang (2021), and Hu et al. (2023), have examined similar spatial interactions in factor models, the SAPT studied in this paper differs from the existing models for several reasons. First, unlike Pesaran and Tosetti (2011), which focuses on spatial autocorrelation in unobserved errors, our model explicitly captures spatial correlations among panel units. Second, the proposed SAPT model functions as a pure spatial arbitrage pricing factor model without lagged or exogenous variables, distinguishing it from the models in Bai and Li (2021) and Yang (2021), which incorporate exogenous features and assume a homogeneous spatial coefficient. This structure presents challenges in identifying suitable instrumental variables for method-of-moments estimation. Third, we consider both observable and latent factor structures. When factors are observable, our model aligns with the spatial asset pricing models of Kou et al. (2018) and Hu et al. (2023) for financial returns. However, when factors are unobservable, which is not considered in Kou et al. (2018) or Hu et al. (2023), our model extends the statistical and econometric factor models by incorporating spatial interaction terms, capturing additional panel information beyond common latent factors. Fourth, our model accommodates panel dimensions that can grow to infinity, differing from the quasi-maximum likelihood estimation (QMLE) framework in Aquaro et al. (2021) and Hu et al. (2023), where the dimension is fixed. This flexibility enables broader applications in high-dimensional settings.

These distinctive features of the proposed SAPT model introduce additional estimation challenges, making conventional spatial econometric methods inadequate. For models with observable factors, the widely used QMLE approach, as discussed in Lee (2004), Yu et al. (2008), and Bai and Li (2021), often encounters computational difficulties due to the large matrix determinants involved in the likelihood function. These challenges become even more pronounced in high-dimensional settings, especially when estimating numerous unit-specific spatial coefficients. In cases with heteroskedastic disturbances, Lin and Lee (2010) demonstrated that the QML estimator for the spatial autoregressive (SAR) model is inconsistent if heteroskedasticity is ignored. To address this problem,

2

they proposed a GMM estimator, which is computationally more efficient than QMLE. However, the SAPT model considered here lacks lagged or exogenous variables, making it difficult to identify suitable instrumental variables for constructing sufficiently many estimating equations.

In view of this, we propose a ridge-regularized Yule-Walker estimator that integrates shrinkage techniques with method-of-moments. By incorporating lagged common factors as instrumental variables, we reformulate parameter estimation as a system of $L_2$-penalty Yule-Walker equations for each panel component, thereby addressing the issue of insufficient number of estimating equations in settings without exogenous variables or structural constraints on spatial effects. In contrast to the regularized method-of-moments approaches proposed by Liao (2013) and Carrasco and Tchuente (2015), which primarily focus on selecting instruments or moment conditions, our method applies ridge regularization directly to the Yule-Walker equations to mitigate potential singularity and improve estimation robustness. We establish the asymptotic properties of our estimator in the setting where both the dimension $N$ and the sample size $T$ approach infinity. Despite the bias inherent in ridge estimators, we demonstrate the feasibility of conducting joint parameter inference. This contrasts with QML estimators in Aquaro et al. (2021) and Hu et al. (2023) which often require finite $N$ and inevitably accumulate asymptotic bias as $N$ diverges with $T$; See Remarks 6 and 7 in Aquaro et al. (2021) for a discussion. While alternative methods, such as those proposed by Bai and Li (2021), ensure parameter consistency under structural constraints and complex bias correction, it remains unclear whether their approach is feasible for the SAPT considered in this paper with heterogeneous spatial interactions.

In the presence of latent factors, our model can be reformulated as an approximate factor model. We propose a two-step procedure to extract latent factors and to estimate unknown parameters. Given the white noise assumption on the error terms in the SAPT model (see Kou et al. (2018), Aquaro et al. (2021), and Hu et al. (2023)), we first apply the auto-covariance-based eigenanalysis approach from Lam and Yao (2012) and Gao and Tsay (2022) to estimate dynamically dependent factors. This ensures that the factors and their lagged counterparts remain uncorrelated with the noise terms, enabling their use as instrumental variables. Once the factors are extracted via eigenanalysis, we implement the Yule-Walker estimation method, replacing the unknown factors with their estimated counterparts. Furthermore, we establish the asymptotic properties of the estimated factors, scalar coefficients, and loading vectors as both the dimension $N$ and sample size $T$ approach infinity. Notably, we also derive the limiting distributions of the estimated factors under a proper rotation matrix, a result not presented in Lam and Yao (2012), offering independent interest for readers.

We conduct extensive simulations to evaluate the accuracy of our estimation method, particularly in estimating the *spatial rho* and the loading matrix, while examining the convergence and asymptotic properties of the jointly estimated parameters. Moreover, we compare our method's predictive performance with QML estimators. The results show that our approach outperforms these alternatives in out-of-sample forecasting. Empirically, we apply our method to two real datasets on U.S. stock returns and housing prices, respectively. In both cases, it achieves superior

out-of-sample forecasting performance compared to QMLE and the classical Fama-French factor model, reinforcing its practical advantages in high-dimensional economic and financial analysis.

This paper makes several significant contributions. First, rather than relying on a mathematical formulation, we derive the SCAPM from a classical mean-variance perspective and extend it to a SAPT framework by integrating a multifactor structure. This approach offers a new perspective for economists and practitioners in understanding spatial asset pricing theory. Second, from a modeling standpoint, the proposed framework is flexible, accommodating both observable and latent factors. This extension provides an opportunity to explore the dynamics of large-dimensional economic and financial panel systems. Third, from a methodological perspective, since QMLE methods require extensive computation and may be impractical in high-dimensional settings with general covariance structures, we propose a shrinkage estimation approach with joint inferential theory for the proposed models. While individual estimators may not be consistent, joint estimation allows for consistent inference. Our procedure is computationally efficient and avoids the need for restrictive distributional or covariance assumptions when using the Yule-Walker estimation method. More importantly, the proposed shrinkage estimation method outperforms the QMLE method in out-of-sample evaluations, highlighting their empirical advantages in high-dimensional applications.

The remainder of the paper is structured as follows. Section 2 outlines the formulations of the SCAPM and SAPT models under study. Section 3 provides the modeling framework and its estimation procedure. Section 4 establishes the asymptotic properties of the derived estimators. Section 5 evaluates the finite-sample performance of the proposed approach through simulations and Section 6 illustrates the proposed model and method with two empirical applications. Section 5 concludes. All proofs and derivations for the asymptotic results are relegated to an online Appendix.

**Notation:** We use the following notation. For a $p \times 1$ vector $\mathbf{u} = (u_1, ..., u_p)'$, $\|\mathbf{u}\|_1 = \sum_{i=1}^{p} |u_i|$ is the $\ell_1$-norm and $\|\mathbf{u}\|_\infty = \max_{1 \le i \le p} |u_i|$ is the $\ell_\infty$-norm. $\mathbf{I}_p$ denotes the $p \times p$ identity matrix. For a matrix $\mathbf{H}$, its Frobenius norm is $\|\mathbf{H}\| = [\text{trace}(\mathbf{H}'\mathbf{H})]^{1/2}$ and its operator norm is $\|\mathbf{H}\|_2 = \sqrt{\lambda_{\max}(\mathbf{H}'\mathbf{H})}$, where $\lambda_{\max}(\cdot)$ denotes the largest eigenvalue of a matrix, and $\|\mathbf{H}\|_{\min}$ is the square root of the minimum non-zero eigenvalue of $\mathbf{HH}'$. $|\mathbf{H}|$ denotes the absolute value of $\mathbf{H}$ elementwisely. The superscript $'$ denotes the transpose of a vector or matrix. We also use the notation $a \asymp b$ to denote $a = O(b)$ and $b = O(a)$ or $a$ and $b$ have the same order of stochastic bound when they are random variables.

## 2 Spatial CAMP and Spatial APT

In this section, we develop a Spatial Capital Asset Pricing Model (SCAPM) using mean-variance analysis within a *complete market* framework. Additionally, we construct a spatial arbitrage pricing theory model by incorporating a multifactor structure.

## 2.1 Complete Market Assumption

We consider a one-period economy with $N$ risky assets in the market whose random returns are denoted as $\mathbf{r} = (r_1, ..., r_N)'$ over the period. The expected return is $\boldsymbol{\mu} = (\mu_1, ..., \mu_N)'$ and the risk-free asset return is $r_f$. Let $r_M$ be the return of the market portfolio or the tangency portfolio in the mean-variance framework of Markowitz (1952) with expected return $\mu_M$.

Suppose the $N$ assets are diverse, with $N$ sufficiently large, and encompass a wide spectrum of asset categories. It is reasonable to assume that each individual asset exhibits some degree of association with the others. Within this framework, we introduce the concept of a *complete market*, which implies that the extensive set of assets enables the formation of suitable linear combinations to replicate the returns of any specific asset in the market. The definition of a *complete market*, or equivalently, a *minimum complete market*, is provided in Definition 1 and equivalently in Definition 2 below.

**Definition 1** (*Complete Market*)**.** *Suppose there are $N$ risky assets in the market, where $N$ is sufficiently large. The market is said to be complete if the return of any asset $r_j$ can be expressed as a linear combination of the remaining $N - 1$ assets, i.e., those indexed by $\{1, ..., N\} \setminus \{j\}$, for $j = 1, ..., N$, However, $r_j$ cannot be replicated using only $N - 2$ assets from $\{1, ..., N\} \setminus \{j\}$.*

**Definition 2** (*Minimum Complete Market*)**.** *A market is called a minimum complete market if it contains at least $N - 1$ assets from a complete market, as the return of the remaining asset can be fully replicated by a linear combination of the other $N - 1$ assets in this minimum complete market.*

This conceptualization of a complete market aligns with the idea that the abundance and diversity of assets enable the construction of portfolios capable of replicating the performance of any individual asset. It suggests that the richness of the market, in terms of asset variety, allows for the creation of synthetic versions of assets by leveraging a diverse set of available instruments.

In a high-dimensional setting, market completeness arises from the vast number and diversity of assets, facilitating the construction of well-diversified portfolios that can closely approximate the returns of specific assets. The concept of a complete market is closely tied to the absence of arbitrage opportunities. In such a market, no-arbitrage conditions ensure that riskless profits cannot be generated through linear combinations of available assets. If arbitrage opportunities existed, they would indicate an incomplete market, as investors could exploit them to create new assets beyond those initially available.

However, achieving a truly complete market in practice is challenging. Real-world markets often face limitations in asset variety, and factors such as transaction costs, market frictions, and short-selling constraints can hinder perfect asset replication. Nonetheless, the notion of a complete market provides a framework for understanding the relationships among assets and their pricing dynamics in a diversified financial environment.
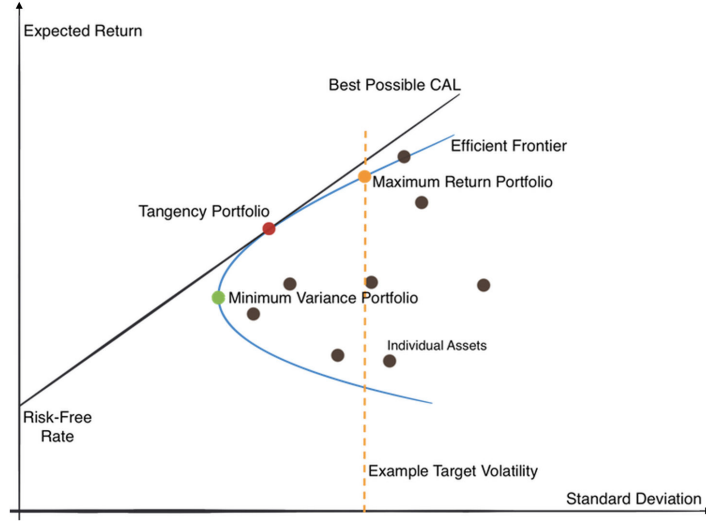
**Figure 1:** Mean-variance efficient frontier with a risk-free asset. The horizontal axis denotes the standard deviation of the portfolio and the vertical axis denotes the expected return of the corresponding portfolio. Available at `https://quantpedia.com/markowitz-model/`.

## 2.2 From CAPM to Spatial CAPM

Based on the mean-variance analysis (e.g., Cochrane (2009)), there exists a weight vector $\boldsymbol{\theta} = (\theta_1, ..., \theta_N)'$ such that the market (or tangency) portfolio can be expressed as $r_M = \boldsymbol{\theta}'\mathbf{r}$, as illustrated in Figure 1. For the $j$-th asset with return $r_j$ and expected return $\mu_j$, the capital asset pricing model (CAPM) of Sharpe (1964) states that

$$\mu_j - r_f = \frac{\text{Cov}(r_j, r_M)}{\text{Var}(r_M)}(\mu_M - r_f),$$

where the quantity $\beta_j = \frac{\text{Cov}(r_{j,t}, r_{M,t})}{\text{Var}(r_{M,t})}$ is referred to as the market beta of the $j$th asset in the finance literature. In practice, the S&P 500 index return often serves as a proxy for the market portfolio, and the market beta can be estimated by running an OLS regression over $T$ periods. For further details, see Chapters 5 and 9 of Cochrane (2009).

Next, we formulate a spatial capital asset pricing model (SCAPM), building on the mean-variance analysis within a complete market defined in Definition 1. For each $j$, we remove $r_j$ from the portfolio return vector $\mathbf{r}$ and consider the mean-variance analysis of the remaining $N-1$ risky assets and the risk-free rate $r_f$. Through the classic mean-variance optimization, we obtain the portfolio weight $\mathbf{w}_j$, where the $j$-th position of $\mathbf{w}_j$ is zero and $\mathbf{w}_j'\mathbf{1}_N = 1$ such that the portfolio $\mathbf{w}_j'\mathbf{r}$ is a tangency portfolio, as illustrated in Figure 1 without the $j$-th asset. The optimal portfolios lie along the capital allocation line (CAL) in the mean-variance framework, with a slope

$$\frac{\mu_{j,M} - r_f}{\sigma_{j,M}},$$

6

where

$$\mu_{j,M} = E(\mathbf{w}'_j\mathbf{r}), \text{ and } \sigma_{j,M} = \sqrt{\text{Var}(\mathbf{w}'_j\mathbf{r})}.$$

Then, for the asset $j$ with expected return $\mu_j$, we have the following theorem.

**Theorem 1.** *Suppose the $N$ risky assets are in a complete market, as described in Definition 1. For the $j$-th risky asset with expected return $\mu_j$, we have the following relationship:*

$$\mu_j - r_f = \frac{\text{Cov}(r_j, \mathbf{w}'_j\mathbf{r})}{\text{Var}(\mathbf{w}'_j\mathbf{r})}(\mu_{j,M} - r_f),$$

*where $\mu_{j,M}$ is the expected return of the tangency portfolio $r_{j,M} = \mathbf{w}'_j\mathbf{r}$ with the $j$-th asset excluded from the portfolio. We define $\rho_j = \frac{\text{Cov}(r_j, \mathbf{w}'_j\mathbf{r})}{\text{Var}(\mathbf{w}'_j\mathbf{r})}$ and refer to it as the "spatial rho" for the $j$-th asset.*

The proof of Theorem 1 can be found in the Appendix. From Theorem 1, we observe that the spatial rho is asset-specific, similar to the market beta in the CAPM. However, the key difference is that the spatial tangent portfolio is also asset-specific, which contrasts with the classical CAPM, where the market portfolio is fixed and unique for all assets.

## 2.3 Spatial Arbitrage Pricing Theory

In this section, we derive a spatial arbitrage pricing theory model following the framework of Ross (1976). To better illustrate the application of the proposed model in asset pricing, we use the notation $\mathbf{r}_t = (r_{1,t}, ..., r_{N,t})'$ to denote a vector of returns to $N$ risky assets at time $t$. Letting $\boldsymbol{\mu}_0 = (\mu_{0,1}, ..., \mu_{0,N})'$ be the expected returns of $\mathbf{r}_t$, we consider the following asset pricing model with spatial interactions and multi-factors:

$$\mathbf{r}_t = \mathbf{D}(\boldsymbol{\rho})\mathbf{W}\mathbf{r}_t + \boldsymbol{\nu}_0 + \mathbf{B}\mathbf{f}_t + \boldsymbol{\varepsilon}_t, \tag{1}$$

where $\mathbf{f}_t = (f_{1,t},...,f_{K,t})'$ consists of $K$ observable factors for which the expected return of $f_{i,t}$ is $\mu_i$, for $i = 1, ..., K$. The columns of $\mathbf{B} = (\boldsymbol{\delta}_1, ..., \boldsymbol{\delta}_K)$ are the associated $K$ loading vectors of the $K$ factors, and $\boldsymbol{\nu}_0 = (\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})\boldsymbol{\mu}_0$, where $\boldsymbol{\rho}$ is a vector of spatial rhos. $\mathbf{W}$ is a known spatial weight matrix with zero main diagonal elements, and $\mathbf{D}(\boldsymbol{\rho}) = diag(\rho_1, ..., \rho_N)$, where $\rho_j$ can be estimated by the method in Section 3 below. We may assume each row of $\mathbf{W}$, denoted as $\mathbf{w}_j$, can be calculated either based on some economic distance or through the mean-variance analysis. We introduce some notations before the derivation of the spatial arbitrage pricing theory. We use $\mathbf{1}_N$ to denote the $N$-dimensional vector of 1, e.g., $\mathbf{1}_N = (1, ..., 1)' \in R^N$. Let $\boldsymbol{\theta} = (\theta_1, ..., \theta_N)'$ represent the weight vector that will be used to construct an arbitrage portfolio. Our derivation proceeds in the following three steps.

*Step 1.* Suppose the random vector of returns $\mathbf{r}_t$ satisfies Model (1). We use a weight vector $\boldsymbol{\theta}$ to construct an arbitrage portfolio of $N$ assets, where we assume $\boldsymbol{\theta}'\mathbf{1}_N = 0$, implying that there is no wealth invested in the portfolio. We also require $\boldsymbol{\theta}$ to be a well-diversified portfolio weight with

each component $\theta_i$ being of order $1/N$ in magnitude as in Ross (1976).

*Step 2.* The random return of the portfolio can be written as

$$\boldsymbol{\theta}'\mathbf{r}_t = \boldsymbol{\theta}'\boldsymbol{\mu}_0 + \boldsymbol{\theta}'(\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1}\mathbf{B}\mathbf{f}_t + \boldsymbol{\theta}'(\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1}\boldsymbol{\varepsilon}_t,$$

where $\boldsymbol{\mu}_0 = (\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1}\boldsymbol{\nu}_0$. We further assume that $\varepsilon_{i,t}$'s are independent with each other, for $i$ and $t$, which is a commonly used assumption in the spatial econometrics literature, and each element of $\mathbf{S}(\boldsymbol{\rho})^{-1} = (\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1}$ are of order $1/N$ in absolute magnitude. Together with Assumption 2 in Section 4 below, by the law of large numbers, we can show that

$$\boldsymbol{\theta}'(\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1}\boldsymbol{\varepsilon}_t = o_p(1),$$

and, hence,

$$\boldsymbol{\theta}'\mathbf{r}_t \approx \boldsymbol{\theta}'\boldsymbol{\mu}_0 + \boldsymbol{\theta}'(\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1}\mathbf{B}\mathbf{f}_t.$$

*Step 3.* If we require that the arbitrage portfolio with weight $\boldsymbol{\theta}$ be chosen with no systematic risk, then

$$\boldsymbol{\theta}'(\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1}\boldsymbol{\delta}_i = 0, \quad i = 1, ..., K. \tag{2}$$

This condition ensures that the return of the arbitrage portfolio becomes $\boldsymbol{\theta}'\boldsymbol{\mu}_0$. Using the constraint of no wealth that $\boldsymbol{\theta}'\mathbf{1}_N = 0$, the return must be zero to prevent arbitrarily large disequilibrium positions. Therefore, we have

$$\boldsymbol{\theta}'\boldsymbol{\mu}_0 = 0. \tag{3}$$

From the relationships in (2), (3), and $\boldsymbol{\theta}'\mathbf{1}_N = 0$, we conclude that $\boldsymbol{\mu}_0$, $\mathbf{1}_N$, and $(\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1}\boldsymbol{\delta}_i$ are on the same hyperplane, for $i = 1, ..., K$. Then there exist $\gamma_0, \gamma_1,..., \gamma_K$ such that

$$\boldsymbol{\mu}_0 = \gamma_{0,i}\mathbf{1}_N + \gamma_i(\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1}\boldsymbol{\delta}_i, \quad i = 1, ..., K. \tag{4}$$

We will solve the above equations for $\gamma_{0,i}$ and $\gamma_i$ by a plug-in method. Note that when $\boldsymbol{\mu}_0 = r_f\mathbf{1}_N$, the return vector of a risk-free asset, the loadings associated with the factors are zero, i.e., $\boldsymbol{\delta}_i = \mathbf{0}$, for $i = 1, ..., K$. Furthermore, if we take $\mathbf{r}_t = f_{i,t}\mathbf{1}_N$, then $\boldsymbol{\mu}_0 = \mu_i\mathbf{1}_N$, and the spatial parameter $\boldsymbol{\rho} = \mathbf{0}$, since there is no spatial effect for a single asset. In this case, the exposure to the $i$-th factor is $\boldsymbol{\delta}_i = \mathbf{1}_N$, while the exposures to the other factors are zero. These special cases result in the following equations:

$$\begin{cases} r_f\mathbf{1}_N = \gamma_{0,i}\mathbf{1}_N, \\ \mu_i\mathbf{1}_N = \gamma_{0,i}\mathbf{1}_N + \gamma_i\mathbf{1}_N, \quad i = 1, ..., K. \end{cases}$$

It follows from the above equations that

$$\gamma_{0,i} = r_f, \ \gamma_i = \mu_i - r_f, \quad i = 1, ..., K,$$

where $\gamma_{0,i}$ turns out to be independent of $i$. Then, (4) becomes

$$\boldsymbol{\mu}_0 = r_f \mathbf{1}_N + (\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1} \boldsymbol{\delta}_1 (\mu_1 - r_f) + ... + (\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1} \boldsymbol{\delta}_K (\mu_K - r_f),$$

or equivalently,

$$(\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})(\boldsymbol{\mu}_0 - r_f \mathbf{1}_N) = \boldsymbol{\delta}_1 (\mu_1 - r_f) + ... + \boldsymbol{\delta}_K (\mu_K - r_f),$$

which is a spatial APT model that extends the SCAMP in Section 2.2 with multi-factors, where $\mu_i - r_f$ is the risk premium of the $i$-th factor and $\boldsymbol{\mu}_0 - r_f \mathbf{1}_N$ is the vector of $N$ excessive asset returns. For the $j$-th asset, we can derive that

$$\mu_{0,j} - r_f = \rho_j \mathbf{w}_j'(\boldsymbol{\mu}_0 - r_f \mathbf{1}_N) + \boldsymbol{\delta}_{1,j}(\mu_1 - r_f) + ... + \boldsymbol{\delta}_{K,j}(\mu_K - r_f), \ \ j = 1, ..., N. \quad (5)$$

Therefore, we may construct a new asset-specific factor, called the spatial factor, defined as $\mathbf{w}_j'(\boldsymbol{\mu}_0 - r_f \mathbf{1}_N)$ associated with the $j$-th asset where the $j$-th element of $\mathbf{w}_j$ is zero according to the definition of the spatial weight. The scalar $\rho_j$ represents the spatial effect on the $j$-th asset, which is termed the *spatial rho*, in contrast to the *market beta* in the classic CAPM of Sharpe (1964).

In the next section, we examine a general APT model that incorporates spatial interactions and propose a Yule-Walker estimation and inference method using factor instruments and ridge techniques for the model.

## 3　General Model and Methodology

### 3.1　Setup

Let $\mathbf{y}_t = (y_{1,t}, ..., y_{N,t})'$ be an $N$-dimensional observable panel of time series at time $t$, where we assume all the data are centered with zero mean. Thus, $\mathbf{y}_t$ replaces $(\mathbf{r}_t - \boldsymbol{\mu}_0)$ in Model (1), and the factors $\{\mathbf{f}_t\}$, for $t = 1, ..., T$, are assumed to have zero mean. Based on the SAPT model in Section 2.3, we assume that $\mathbf{y}_t$ follows the following general structure:

$$\mathbf{y}_t = \mathbf{D}(\boldsymbol{\rho})\mathbf{W}\mathbf{y}_t + \mathbf{B}\mathbf{f}_t + \boldsymbol{\varepsilon}_t, \ \ t = 1, ..., T, \quad (6)$$

where $\mathbf{f}_t$ is a $K$-dimensional factor process that is either observable or unobservable, $\mathbf{B}$ is the loading matrix associated with the factors, $\mathbf{W}$ is the $N \times N$ spatial weight matrix that measures the dependence among different economic units or individuals of $\mathbf{y}_t$. $\mathbf{D}(\boldsymbol{\rho}) = \mathrm{diag}(\rho_1, ..., \rho_N)$, where $\rho_j$ is an unknown coefficient parameter for the $j$-th individual. $\boldsymbol{\varepsilon}_t$ is a white noise term that is uncorrelated with $\mathbf{f}_t$, but we allow for dependence between $\mathbf{f}_{t+j}$ and $\boldsymbol{\varepsilon}_t$, for $j \geq 1$, since the factors $\mathbf{f}_t$'s are usually serially dependent, which may be correlated with some lagged noise terms.

It is a common practice in spatial econometrics to assume that $\mathbf{W}$ is known, and the main diagonal elements of $\mathbf{W}$ are zero. The weights may be based on physical distance, social networks,

or "economic" distance, as seen in Case et al. (1993). For example, we may take $w_{ij} := (s_i d_{ij})^{-1}$, for $i \neq j$, and $w_{ii} = 0$, where $d_{ij}$ is the physical distance between location $i$ and location $j$, and $s_i := \sum_j d_{ij}^{-1}$. Alternatively, we may take $d_{ij}^{-1}$ as the sample correlation between the $i$-th and $j$-th economic units when there is no clear physical distance between them. When $\rho_1 = \ldots = \rho_N$, the spatial interaction term in Model (6) reduces to the classical setting in the spatial econometrics literature, such as Lee (2004), among others.

For a given spatial weight matrix $\mathbf{W} = (\mathbf{w}_1, \ldots, \mathbf{w}_N)'$, where $\mathbf{w}_i$ is the $i$-th row vector of $\mathbf{W}$, our goal is to estimate the unknown coefficients in $\boldsymbol{\rho}$ and $\mathbf{B}$ when the factors $\mathbf{f}_t$'s are observable. When the factors $\mathbf{f}_t$'s are latent, we also need to recover the latent factors.

## 3.2 Shrinkage Yule-Walker Estimation with Observed Factors

In this section, we study the scenario when the factors are observed and propose a generalized shrinkage Yule-Walker method to estimate the unknown coefficients, which is essentially a combination of ridge regression and the method-of-moments. To this end, we begin with some useful notation. Define $\boldsymbol{\Sigma}_{yf}(k) = \text{Cov}(\mathbf{y}_t, \mathbf{f}_{t-k})$ as the covariance matrix between $\mathbf{y}_t$ and the past lagged factor variables $\mathbf{f}_{t-k}$, and $\boldsymbol{\Sigma}_f(k) = \text{Cov}(\mathbf{f}_t, \mathbf{f}_{t-k})$ as the lag-$k$ auto-covariance matrix of $\mathbf{f}_t$, for $k \geq 0$. Then, Model (6) implies that

$$\boldsymbol{\Sigma}_{yf}(k) = \mathbf{D}(\boldsymbol{\rho})\mathbf{W}\boldsymbol{\Sigma}_{yf}(k) + \mathbf{B}\boldsymbol{\Sigma}_f(k), \quad k \geq 0. \tag{7}$$

Let $\mathbf{e}_i$ be the $i$th unit vector with the $i$th element equal to 1 and other elements being zero. For each $k \geq 0$, it follows from (7) that

$$\mathbf{e}_i'\boldsymbol{\Sigma}_{yf}(k) = \mathbf{e}_i'\mathbf{D}(\boldsymbol{\rho})\mathbf{W}\boldsymbol{\Sigma}_{yf}(k) + \mathbf{e}_i'\mathbf{B}\boldsymbol{\Sigma}_f(k), \quad i = 1, \ldots, N. \tag{8}$$

Note that $\mathbf{e}_i'\mathbf{D}(\boldsymbol{\rho})\mathbf{W} = \rho_i \mathbf{w}_i'$ and $\mathbf{e}_i'\mathbf{B} = \mathbf{b}_i'$, where $\mathbf{w}_i$ and $\mathbf{b}_i$ are the $i$th row vectors of $\mathbf{W}$ and $\mathbf{B}$, respectively. Then, (8) becomes

$$\boldsymbol{\Sigma}_{yf}'(k)\mathbf{e}_i = \boldsymbol{\Sigma}_{yf}'(k)\mathbf{w}_i \rho_i + \boldsymbol{\Sigma}_f'(k)\mathbf{b}_i, \quad i = 1, \ldots, N. \tag{9}$$

In practice, given the sample data $\{(\mathbf{y}_t, \mathbf{f}_t) : t = 1, \ldots, T\}$, by a similar argument to the Yule-Walker estimation method with a given lag $k \geq 0$, we may solve the following minimization problem:

$$(\widehat{\rho}_i, \widehat{\mathbf{b}}_i')' = \arg \min_{\rho \in R, \mathbf{b} \in R^r} \{\|\widehat{\boldsymbol{\Sigma}}_{yf}'(k)\mathbf{e}_i - \widehat{\boldsymbol{\Sigma}}_{yf}'(k)\mathbf{w}_i \rho - \widehat{\boldsymbol{\Sigma}}_f'(k)\mathbf{b}\|_2^2\}, \ i = 1, \ldots, N, \tag{10}$$

where

$$\widehat{\boldsymbol{\Sigma}}_{yf}(k) = \frac{1}{T} \sum_{t=k+1}^{T} \mathbf{y}_t \mathbf{f}_{t-k} \quad \text{and} \quad \widehat{\boldsymbol{\Sigma}}_f(k) = \frac{1}{T} \sum_{t=k+1}^{T} \mathbf{f}_t \mathbf{f}_{t-k}'$$

are the sample versions of $\boldsymbol{\Sigma}_{yf}(k)$ and $\boldsymbol{\Sigma}_f(k)$, respectively. For each $i$, we observe that there are $K + 1$ unknown coefficients in the optimization problem (10), but there are only $K$ equations for

each lag $k$ in (10), which implies that the optimization problem is not well-defined if we only make use of a single $k$ in the Yule-Walker estimation. To see this, we cast problem (10) into the framework of the generalized method of moments (GMM) (Hansen (1982)). Let $\mathbf{f}_{t-k}$ be the instrument, the moment conditions for (6) are

$$E\mathbf{h}_{t,k}(\rho_i, \mathbf{b}_i) = 0, \text{ where } \mathbf{h}_{t,k}(\rho_i, \mathbf{b}_i) = (y_{i,t} - \rho_i\mathbf{w}_i'\mathbf{y}_t - \mathbf{b}_i'\mathbf{f}_t)\mathbf{f}_{t-k}', \ i = 1, ..., N,$$

which is equivalent to that in (9) for each $k$. When $k = 0$, it is not hard to see that

$$\mathbf{h}_{t,0}(\rho_i, \mathbf{b}_i) = \frac{\partial \varepsilon_{i,t}(\rho_i, \mathbf{b}_i)}{\partial \mathbf{b}_i},$$

where $\varepsilon_{i,t}(\rho_i, \mathbf{b}_i) = (y_{i,t} - \rho_i\mathbf{w}_i'\mathbf{y}_t - \mathbf{b}_i'\mathbf{f}_t)^2$. Therefore, the equations produced by only taking partial derivatives concerning parameter $\mathbf{b}_i$ are not sufficient to estimate $\rho_i$ and $\mathbf{b}_i$ simultaneously. Then we conclude that estimation equations in (10) are not sufficient for any given $k \geq 0$.

To address this challenge, we combine two sets of estimating equations, resulting in $2K$ equations, which exceed the number of parameters by $K + 1$ when $K \geq 1$. We first note the importance of cross-sectional dependence in asset returns and economic data, so we retain the $k = 0$ equations, which align with GMM using $\mathbf{f}_t$ as instruments and the Least-Squares method. Additionally, since short-term dependence is more significant than long-term in economic and financial data, we focus on Equation (10) with $k = 0$ and $k = 1$. These lags capture the key dynamic information, reflecting the most relevant dependencies while excluding less impactful higher lags.

Specifically, let

$$\widehat{\mathbf{Y}}_i = \begin{pmatrix} \widehat{\boldsymbol{\Sigma}}_{yf}'\mathbf{e}_i \\ \widehat{\boldsymbol{\Sigma}}_{yf}'(1)\mathbf{e}_i \end{pmatrix} \text{ and } \widehat{\mathbf{X}}_i = \begin{pmatrix} \widehat{\boldsymbol{\Sigma}}_{yf}'\mathbf{w}_i & \widehat{\boldsymbol{\Sigma}}_f \\ \widehat{\boldsymbol{\Sigma}}_{yf}'(1)\mathbf{w}_i & \widehat{\boldsymbol{\Sigma}}_f'(1) \end{pmatrix},$$

where $\widehat{\boldsymbol{\Sigma}}_{yf} = \widehat{\boldsymbol{\Sigma}}_{yf}(0)$ and $\widehat{\boldsymbol{\Sigma}}_f = \widehat{\boldsymbol{\Sigma}}_f(0)$. Due to the spatial nature of Model (6), $\widehat{\mathbf{X}}_i$ is asymptotically singular, though not in finite samples. To address this, we apply ridge regression. Define $\boldsymbol{\beta} = (\widehat{\rho}, \mathbf{b}')' \in \mathbb{R}^{K+1}$ and solve the following optimization problem for a given $\lambda_i > 0$:

$$\widehat{\boldsymbol{\beta}}_i(\lambda_i) = (\widehat{\rho}_i, \widehat{\mathbf{b}}_i')' = \arg\min_{\rho \in R, \mathbf{b} \in R^K}\{\|\widehat{\mathbf{Y}}_i - \widehat{\mathbf{X}}_i\boldsymbol{\beta}\|_2^2 + \lambda_i\|\boldsymbol{\beta}\|_2^2\}, \ i = 1, ..., N. \tag{11}$$

The estimator has the explicit form:

$$\widehat{\boldsymbol{\beta}}_i(\lambda_i) = (\widehat{\mathbf{X}}_i'\widehat{\mathbf{X}}_i + \lambda_i\mathbf{I}_{K+1})^{-1}\widehat{\mathbf{X}}_i'\widehat{\mathbf{Y}}_i, \ i = 1, ..., N, \tag{12}$$

which is the ridge estimator. In the subsequent analysis, we denote $(\widehat{\mathbf{X}}_i'\widehat{\mathbf{X}}_i)^+$ as the Moore-Penrose generalized inverse and let $\widehat{\boldsymbol{\beta}}_i = \widehat{\boldsymbol{\beta}}_i(0)$ as $\lambda_i \to 0$. Theorem 3 establishes the joint asymptotic distribution under these conditions, enabling joint inference. In finite samples, the estimator depends on the number of lagged auto-covariances used in the Yule-Walker estimation in (10), but its asymptotic convergence remains valid, as demonstrated in Section 4.

## 3.3 Boosting the Strength of Factor Instruments

In practice, we stack only the cases when $k = 0$ and $k = 1$ as discussed in Section 3.2. This approach provides a jointly consistent estimator and avoids unnecessary errors in the generalized method of moments estimation. The choice of $k = 1$ is based on the assumption that short-term dependence is often more relevant than long-term dependence. If necessary, we can define a measure to select the optimal lag $k^*$ as follows:

$$k^* = \arg \max_{1 \leq k \leq \bar{k}} |\det(\widehat{\mathbf{\Sigma}}_f(k))|,$$

where $\bar{k}$ is a small positive integer, and $|\det(\widehat{\mathbf{\Sigma}}_f(k))|$ is the product of the singular values of $\widehat{\mathbf{\Sigma}}_f(k)$. This measure captures the correlation strength between the lagged instruments and the contemporaneous factors. We then define

$$\widehat{\mathbf{Y}}_{i,*} = \begin{pmatrix} \widehat{\mathbf{\Sigma}}'_{yf} \mathbf{e}_i \\ \widehat{\mathbf{\Sigma}}'_{yf}(k^*) \mathbf{e}_i \end{pmatrix} \text{ and } \widehat{\mathbf{X}}_{i,*} = \begin{pmatrix} \widehat{\mathbf{\Sigma}}'_{yf} \mathbf{w}_i & \widehat{\mathbf{\Sigma}}_f \\ \widehat{\mathbf{\Sigma}}'_{yf}(k^*) \mathbf{w}_i & \widehat{\mathbf{\Sigma}}'_f(k^*) \end{pmatrix},$$

and substitute them into (11), yielding the refined estimator

$$\widehat{\boldsymbol{\beta}}_{i,*}(\lambda_i) = (\widehat{\mathbf{X}}'_{i,*} \widehat{\mathbf{X}}_{i,*} + \lambda_i \mathbf{I}_{K+1})^{-1} \widehat{\mathbf{X}}'_{i,*} \widehat{\mathbf{Y}}_{i,*}, \quad i = 1, ..., N. \tag{13}$$

## 3.4 Estimation When Factors Are Latent

In this section, we address the case where the factor processes $\mathbf{f}_t$ are unobservable. We focus on latent factors that represent the internal dynamics driving the data $\mathbf{y}_t$, because the case with factors arising from some external data sources is similar to the diffusion-index framework in Stock and Watson (2002) and Gao and Tsay (2024). The estimation of the proposed model becomes more complex because, in addition to estimating the parameters in $\boldsymbol{\rho}$ and $\mathbf{B}$, we must also recover the unknown factors. Note that

$$\mathbf{y}_t = (\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1} \mathbf{B} \mathbf{f}_t + (\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1} \boldsymbol{\varepsilon}_t = \mathbf{\Lambda} \mathbf{f}_t + \boldsymbol{\xi}_t, \tag{14}$$

where $\mathbf{\Lambda} = (\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1} \mathbf{B}$ and $\boldsymbol{\xi}_t = (\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1} \boldsymbol{\varepsilon}_t$ are the loading matrix associated with the factor process $\mathbf{f}_t$ and the idiosyncratic term, respectively. Symbolically, (14) is a factor model with unknown factors and loading matrix, both of which need to be estimated from the data $\mathbf{y}_t$, for $t = 1, ..., T$.

Under the framework of Model (14), we have a factor model with static factors and could use either the PCA method of Bai and Ng (2002) to estimate cross-sectional factors or the eigen-analysis method in Lam and Yao (2012) to extract dynamically dependent factors. However, the PCA method is not suitable for the spatial interactions in Model (14), as the idiosyncratic noise recovered by PCA is often serially correlated. In contrast, the noise term $\boldsymbol{\varepsilon}_t$ (or $\boldsymbol{\xi}_t$) in the spatial model is white noise, which contradicts the PCA framework.

For spatial panel dynamic models in econometrics, the noise term $\boldsymbol{\varepsilon}_t$ (or $\boldsymbol{\xi}_t$) is white with zero

serial correlation, while the dynamically dependent factors $\mathbf{f}_t$ capture all the dynamic information of the data $\mathbf{y}_t$. This framework aligns with Lam et al. (2011), Lam and Yao (2012), and Gao and Tsay (2022), among others. Based on the auto-covariance-based eigenanalysis in Lam et al. (2011), we propose a two-step procedure to estimate the factors and other unknown coefficients, assuming the number of factors $K$ is known. The method for determining $K$ will be discussed later.

Note that $\mathbf{\Lambda}$ and $\mathbf{f}_t$ are not uniquely determined in (14) and they require certain identification conditions. For simplicity, we assume that $\mathbf{\Lambda}$ is a semi-orthogonal matrix scaled by $\sqrt{N}$ such that $\mathbf{\Lambda}'\mathbf{\Lambda}/N = \mathbf{I}_K$. However, the loading and factors are still not uniquely identified because we can replace $(\mathbf{\Lambda}, \mathbf{f}_t)$ with $(\mathbf{\Lambda H}, \mathbf{H}'\mathbf{f}_t)$ for any orthonormal matrix $\mathbf{H} \in \mathbb{R}^{K \times K}$. Nevertheless, the linear space spanned by the columns of $\mathbf{\Lambda}$, denoted $\mathcal{M}(\mathbf{\Lambda})$, is uniquely defined and referred to as the factor loading space.

Under the assumption that $\boldsymbol{\varepsilon}_t$ is a white noise process and $\mathrm{Cov}(\mathbf{f}_t, \boldsymbol{\varepsilon}_{t+j}) = 0$, for $j \geq 0$, we allow for the possibility that $\mathbf{f}_t$ may depend on the past lagged noises $\boldsymbol{\varepsilon}_{t-k}$, for some $k \geq 1$, as $\mathbf{f}_t$ is a dynamically dependent process. For any integer $k \geq 1$, define the following covariance matrices of interest:

$$\mathbf{\Sigma}_y(k) = \mathrm{Cov}(\mathbf{y}_t, \mathbf{y}_{t-k}), \ \ \mathbf{\Sigma}_f(k) = \mathrm{Cov}(\mathbf{f}_t, \mathbf{f}_{t-k}), \ \ \text{and} \ \ \mathbf{\Sigma}_{f\xi}(k) = \mathrm{Cov}(\mathbf{f}_t, \boldsymbol{\xi}_{t-k}).$$

From (14), we have

$$\mathbf{\Sigma}_y(k) = \mathbf{\Lambda}\mathbf{\Sigma}_f(k)\mathbf{\Lambda}' + \mathbf{\Lambda}\mathbf{\Sigma}_{f\xi}(k), \quad k \geq 1. \tag{15}$$

For a pre-specified integer $k_0 > 0$, define

$$\mathbf{M} = \sum_{k=1}^{k_0} \mathbf{\Sigma}_y(k)\mathbf{\Sigma}_y'(k) = \mathbf{\Lambda} \sum_{k=1}^{k_0} [\mathbf{\Sigma}_f(k)\mathbf{\Lambda}' + \mathbf{\Sigma}_{f\xi}(k)][\mathbf{\Lambda}\mathbf{\Sigma}_f'(k) + \mathbf{\Sigma}_{f\xi}'(k)]\mathbf{\Lambda}', \tag{16}$$

which is an $N \times N$ semi-positive definite matrix. Let $\mathbf{\Lambda}_c$ denote the orthogonal complement matrix of $\mathbf{\Lambda}$. We observe that $\mathbf{M}\mathbf{\Lambda}_c = \mathbf{0}$, implying that the columns of $\mathbf{\Lambda}_c$ are the eigenvectors corresponding to the zero eigenvalues of $\mathbf{M}$. The factor loading space $\mathcal{M}(\mathbf{\Lambda})$ is thus spanned by the eigenvectors (scaled by $\sqrt{N}$) corresponding to the $K$ non-zero eigenvalues of $\mathbf{M}$. The integer $k_0$ in (16) is a prescribed value that allows us to accumulate dynamic information across different lags. Since the dynamic dependence between $\mathbf{y}_t$ and $\mathbf{y}_{t-k}$ typically decreases as $k$ increases for stationary processes, a small $k_0$ is generally sufficient in practice. For further details on the rationale for using (16) to estimate the loading space from a projection perspective, we refer readers to Gao and Tsay (2021).

In practice, given the sample data $\{\mathbf{y}_t \mid t = 0, 1, \ldots, T\}$, the first step of the procedure is to estimate the loading matrix $\mathbf{\Lambda}$ or its column space $\mathcal{M}(\mathbf{\Lambda})$, and to recover the factor process $\mathbf{f}_t$, assuming that the number of factors $K$ is known. The estimation of $K$ will be discussed later. Let $\widehat{\mathbf{\Sigma}}_y(k)$ denote the lag-$k$ sample autocovariance matrix of $\mathbf{y}_t$, defined similarly to those in (10). To

estimate $\mathcal{M}(\mathbf{\Lambda})$, we perform an eigen-analysis of the sample version of $\mathbf{M}$, defined as

$$\widehat{\mathbf{M}} = \sum_{k=1}^{k_0} \widehat{\mathbf{\Sigma}}_y(k)\widehat{\mathbf{\Sigma}}_y'(k). \tag{17}$$

Let $\widehat{\mathbf{\Lambda}}$ be the standardized semi-orthogonal matrix consisting of the eigenvectors of $\widehat{\mathbf{M}}$, scaled by $\sqrt{N}$, as its columns. The recovered factor processes are denoted as $\widehat{\mathbf{f}}_t = \frac{1}{N}\widehat{\mathbf{\Lambda}}'\mathbf{y}_t$, which can be obtained by the Ordinary Least Squares (OLS) method.

In the second step, we estimate the scalar coefficient vector $\boldsymbol{\rho}$ and the loading matrix $\mathbf{B}$ in Model (6). Let $\widehat{\mathbf{f}}_1, \ldots, \widehat{\mathbf{f}}_T$ denote the estimated factors obtained in the first step. Define the following quantities:

$$\widetilde{\mathbf{\Sigma}}_{yf}(k) = \frac{1}{T}\sum_{t=k+1}^{T}\mathbf{y}_t\widehat{\mathbf{f}}_{t-k}', \ \ \widetilde{\mathbf{\Sigma}}_f(k) = \sum_{t=k+1}^{T}\widehat{\mathbf{f}}_t\widehat{\mathbf{f}}_{t-k}', \ \text{and} \ \widetilde{\mathbf{\Sigma}}_{\varepsilon f}(k) = \frac{1}{T}\sum_{t=k+1}^{T}\boldsymbol{\varepsilon}_t\widehat{\mathbf{f}}_{t-k}',$$

and $\widetilde{\mathbf{\Sigma}}_{yf} = \widetilde{\mathbf{\Sigma}}_{yf}(0)$, $\widetilde{\mathbf{\Sigma}}_f = \widetilde{\mathbf{\Sigma}}_f(0)$, and $\widetilde{\mathbf{\Sigma}}_{\varepsilon f} = \widetilde{\mathbf{\Sigma}}_{\varepsilon f}(0)$. Following a similar procedure to the shrinkage Yule-Walker estimation in Section 3.2, where the factors are observable, we formulate the following optimization problem for the case of augmenting only $k = 0$ and $k = 1$, with a given $\lambda_i > 0$:

$$\widetilde{\boldsymbol{\beta}}_i(\lambda_i) = (\widetilde{\rho}_i, \widetilde{\mathbf{b}}_i')' = \arg\min_{\rho \in R, \mathbf{b} \in R^r}\{\|\widetilde{\mathbf{Y}}_i - \widetilde{\mathbf{X}}_i\boldsymbol{\beta}\|_2^2 + \lambda_i\|\boldsymbol{\beta}\|_2^2\}, \ i = 1, ..., N, \tag{18}$$

where

$$\widetilde{\mathbf{Y}}_i = \begin{pmatrix} \widetilde{\mathbf{\Sigma}}_{yf}'\mathbf{e}_i \\ \widetilde{\mathbf{\Sigma}}_{yf}'(1)\mathbf{e}_i \end{pmatrix} \ \text{and} \ \widetilde{\mathbf{X}}_i = \begin{pmatrix} \widetilde{\mathbf{\Sigma}}_{yf}'\mathbf{w}_i & \widetilde{\mathbf{\Sigma}}_f' \\ \widetilde{\mathbf{\Sigma}}_{yf}'(1)\mathbf{w}_i & \widetilde{\mathbf{\Sigma}}_f'(1) \end{pmatrix}$$

represent the response variables and covariates, respectively. The Yule-Walker estimation in (18) then yields the least squares (LS) estimator for $\boldsymbol{\beta}$ as

$$\widetilde{\boldsymbol{\beta}}_i(\lambda_i) = (\widetilde{\mathbf{X}}_i'\widetilde{\mathbf{X}}_i + \lambda_i\mathbf{I}_{K+1})^{-1}\widetilde{\mathbf{X}}_i'\widetilde{\mathbf{Y}}_i, \ i = 1, ..., N. \tag{19}$$

Thus, we perform $N$ Yule-Walker estimation procedures, for $i = 1, \ldots, N$ and obtain the estimators $\widetilde{\boldsymbol{\rho}} = (\widetilde{\rho}_1, \ldots, \widetilde{\rho}_N)'$ and $\widetilde{\mathbf{B}} = (\widetilde{\mathbf{b}}_1, \ldots, \widetilde{\mathbf{b}}_N)'$. We can similarly define $\widetilde{\boldsymbol{\beta}}_i = \widetilde{\beta}_i(0)$ by adopting the Moore-Penrose inverse of $\widetilde{\mathbf{X}}_i'\widetilde{\mathbf{X}}_i$ as in (12). Theorem 6 in Section 4 establishes the joint asymptotic distribution, which can be utilized for joint inference under this condition.

In practice, we may also use the boosting method described in Section 3.3 to select the optimal lag $k*$. We can then replace $\widetilde{\mathbf{\Sigma}}_{yf}(1)$ and $\widetilde{\mathbf{\Sigma}}_f(1)$ with $\widetilde{\mathbf{\Sigma}}_{yf}(k^*)$ and $\widetilde{\mathbf{\Sigma}}_f(k^*)$, respectively, in $\widetilde{\mathbf{Y}}_i$ and $\widetilde{\mathbf{X}}_i$. The estimator $\widetilde{\boldsymbol{\beta}}_{i,*}(\lambda_i)$ can be obtained in the same manner as that described in Section 3.3.

## 3.5 Selecting the Number of Factors and the Penalty Parameters

In this section, we discuss the determination of the number of factors $K$ in Model (14), which is typically unknown in practice. Over the past decades, several methods have been developed to

estimate $K$, including the information criteria proposed by Bai and Ng (2002), the random matrix theory approach in Onatski (2010), the ratio-based method in Lam and Yao (2012) and Ahn and Horenstein (2013), the canonical correlation analysis technique in Gao and Tsay (2019), and the white noise testing approach in Gao and Tsay (2022), among others. In this paper, we introduce two widely used methods for estimating $K$.

The first method is an information criterion introduced by Bai and Ng (2002). It estimates $K$ by

$$\widehat{K} = \arg\min_{0 \leq j \leq J} \log\left(\frac{1}{NT} \sum_{t=1}^{T} \|\mathbf{y}_t - \frac{1}{N}\widehat{\mathbf{\Lambda}}_j\widehat{\mathbf{\Lambda}}'_j\mathbf{y}_t\|_2^2\right) + jg(T, N), \tag{20}$$

where $J$ is a prescribed upper bound, $\widehat{\mathbf{\Lambda}}_j$ is a $N \times j$ estimated loading matrix, and $g(T, N)$ is a penalty function of $(N, T)$ such that $g(T, N) = o(1)$ and $\min\{N, T\}g(T, N) \to \infty$. Two examples of $g(T, N)$ suggested by Bai and Ng (2002) are IC1 and IC2 given below:

$$IC1 = \frac{N + T}{NT} \log\left(\frac{NT}{N + T}\right) \quad \text{and} \quad IC2 = \frac{N + T}{NT} \log(\min\{N, T\}).$$

For the estimation of $K$, in addition to the information criterion in (20), we can adopt the ratio-based method proposed in Lam and Yao (2012) and Ahn and Horenstein (2013). Let $\widehat{\mu}_1 \geq \cdots \geq \widehat{\mu}_N$ be the $N$ eigenvalues of $\widehat{\mathbf{M}}$. We estimate $K$ by

$$\widehat{K} = \arg\min_{1 \leq l \leq R} \widehat{\mu}_{l+1}/\widehat{\mu}_l, \tag{21}$$

where $R = \lfloor N/2 \rfloor$ is commonly used, as suggested by Lam and Yao (2012).

For the selection of the penalty parameter $\lambda_i$, it is common to assume that $\lambda_i \in \mathcal{S}$, where $\mathcal{S}$ is a candidate set consisting of possible penalty choices. We split the data sample into two segments, $\mathbf{y}_1, \ldots, \mathbf{y}_{T_1}$ and $\mathbf{y}_{T_1+1}, \ldots, \mathbf{y}_T$. Suppose $\widehat{\rho}_i(\lambda)$ and $\widehat{\mathbf{b}}_i(\lambda)$ are the estimators obtained from the first segment. The optimal $\lambda$ is chosen by solving

$$\widehat{\lambda}_i = \text{argmin}_{\lambda \in \mathcal{S}} \frac{1}{T - T_1} \sum_{t=T_1+1}^{T} \|y_{i,t} - \widehat{\rho}_i(\lambda)\mathbf{w}'_i\mathbf{y}_t - \widehat{\mathbf{b}}_i(\lambda)\mathbf{f}_t\|_2^2. \tag{22}$$

When the factors are unobservable, we replace $\mathbf{f}_t$ with $\widehat{\mathbf{f}}_t$, which is estimated from the second segment using the estimator $\widehat{\mathbf{\Lambda}}(\lambda)$ obtained from the first segment.

## 4 Theoretical Properties

In this section, we present the asymptotic theory for the estimation method of Section 3, when both the dimension $N$ and the sample size $T$ tend to infinity. We focus on the estimating equations with lags $k = 0$ and $k = 1$, which typically capture the majority of the cross-sectional and dynamic dependencies in the data. A constant $C$ is used generically, with its value potentially varying across different parts of the analysis. We begin with some assumptions.

**Assumption 1.** *The process $\{(\mathbf{y}_t, \mathbf{f}_t)\}$ is strictly stationary and $\alpha$-mixing with the mixing coefficient satisfying the condition $\sum_{k=1}^{\infty} \alpha_N(k)^{1-2/\gamma} < \infty$ for some $\gamma > 2$, where*

$$\alpha_N(k) = \sup_i \sup_{A \in \mathcal{F}_{-\infty}^i, B \in \mathcal{F}_{i+k}^{\infty}} |P(A \cap B) - P(A)P(B)|,$$

*and $\mathcal{F}_i^j$ is the $\sigma$-field generated by $\{(\mathbf{y}_t, \mathbf{f}_t) : i \leq t \leq j\}$.*

**Assumption 2.** *The spatial weight matrix $\mathbf{W}$ is known with zero main diagonal elements, and the matrix $\mathbf{S}_N(\boldsymbol{\rho}) := \mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W}$ is invertible. The row and column sums of $|\mathbf{W}|$ and $|\mathbf{S}_N(\boldsymbol{\rho})^{-1}|$ are bounded uniformly in $N$.*

**Assumption 3.** *$\{\boldsymbol{\varepsilon}_t\}$ is a white noise process satisfying $\mathrm{Cov}(\mathbf{y}_{t-j}, \boldsymbol{\varepsilon}_t) = \mathbf{0}$ and $\mathrm{Cov}(\mathbf{f}_{t-k}, \boldsymbol{\varepsilon}_t) = \mathbf{0}$, for $j \geq 1$ and $k \geq 0$, respectively.*

**Assumption 4.** *(i) If $\mathbf{f}_t$'s are observed, each element in $\mathbf{B}$ are bounded uniformly in $N$; (ii) If $\mathbf{f}_t$'s are latent, the loading matrix $\mathbf{B}$ is of full rank such that $\frac{1}{N}\mathbf{B}'\mathbf{S}_N'(\boldsymbol{\rho})^{-1}\mathbf{S}_N(\boldsymbol{\rho})^{-1}\mathbf{B} = \mathbf{I}_r$, which is an identity matrix.*

**Assumption 5.** *For $1 \leq j \leq K$ and $1 \leq k \leq N$, $E|f_{j,t}|^{2\gamma} < C$ and $E|\varepsilon_{k,t}|^{2\gamma} < C$, where $\gamma$ is given in Assumption 1.*

**Assumption 6.** *For $i = 1, ..., N$, the rank of matrix $\mathbf{X}_i'\mathbf{X}_i + \lambda\mathbf{I}_{K+1}$ is $K + 1$, for any $\lambda > 0$, where*

$$\mathbf{X}_i = \begin{pmatrix} \boldsymbol{\Sigma}_{yf}'\mathbf{w}_i & \boldsymbol{\Sigma}_f' \\ \boldsymbol{\Sigma}_{yf}'(1)\mathbf{w}_i & \boldsymbol{\Sigma}_f'(1) \end{pmatrix}.$$

Assumption 1 is standard for dependent random processes. See Gao et al. (2019) for a theoretical justification for VAR models. In fact, the assumption of strict stationarity can be removed and we only need to replace definitions of $\boldsymbol{\Sigma}_y(k)$ and $\boldsymbol{\Sigma}_f(k)$ with $\frac{1}{T}\sum_{t=k+1}^T \mathrm{Cov}(\mathbf{y}_t, \mathbf{f}_{t-k})$ and $\frac{1}{T}\sum_{t=k+1}^T \mathrm{Cov}(\mathbf{f}_t, \mathbf{f}_{t-k})$, respectively, and the results still hold throughout the paper. Assumption 2 is commonly used in the spatial econometrics literature to limit the dependence across different locations or economic units; see, for example, Lee and Yu (2010). Assumption 3 is weaker than the independence assumptions imposed in the spatial econometrics literature and we also allow for possible dependence between $\mathbf{y}_{t+j}$ and $\mathbf{f}_{t+k}$ and past lagged of noises, for $j \geq 0$ and $k \geq 1$. Assumption 4 is standard for the loading matrix under the scenarios when the factors are either observed or latent. Assumption 5 imposes some moment conditions on the factors and noise terms. It is not hard to see that $E|y_{i,t}|^{2\gamma} < C$ under Assumptions 2, 4 and 5. Furthermore, this also implies that $E|\mathbf{w}_i'\boldsymbol{\Sigma}_{yf}\mathbf{f}_t|^{2\gamma} < C$, $E|\mathbf{w}_i'\boldsymbol{\Sigma}_{yf}(1)\mathbf{f}_{t-1}|^{2\gamma} < C$, $E\|\boldsymbol{\Sigma}_f\mathbf{f}_t\|_2^{2\gamma} < C$, and $E\|\boldsymbol{\Sigma}_f(1)\mathbf{f}_{t-1}\|_2^{2\gamma} < C$, which are used to establish the convergence of the variance of $\mathbf{S}_{N,T}$, as defined in (IA.5) of the online Appendix. Assumption 6 ensures that the ridge solutions in (12) and (19) are well-defined.

Now, we present the asymptotic properties of $\widehat{\boldsymbol{\beta}}_i$, for $i = 1, \ldots, N$.

**Theorem 2.** *Let Assumptions 1 − 6 hold.*
*(i) If $N = o(T)$, we have*

$$\|\widehat{\boldsymbol{\beta}}_i(\lambda_i) - \widehat{\mathbf{X}}_i(\lambda_i)^{-1}\widehat{\mathbf{X}}_i\widehat{\mathbf{X}}_i'\boldsymbol{\beta}_i\|_2 = O_p(T^{-1/2}), \quad i = 1, ..., N,$$

*as $N, T \to \infty$, where $\widehat{\mathbf{X}}_i(\lambda_i) = \widehat{\mathbf{X}}_i'\widehat{\mathbf{X}}_i + \lambda_i\mathbf{I}_{K+1}$.*
*(ii) If $N = o(T)$ and let $\lambda_i \to 0$, we have*

$$\|(\widehat{\mathbf{X}}_i'\widehat{\mathbf{X}}_i)(\widehat{\boldsymbol{\beta}}_i - \boldsymbol{\beta}_i)\|_2 = O_p(T^{-1/2}), \quad i = 1, ..., N,$$

*as $N, T \to \infty$.*

Theorem 2 implies that the ridge estimator for $\boldsymbol{\beta}$ is biased, which is a common issue in ridge estimation. However, we can establish the joint convergence of $\widehat{\rho}_i$ and $\mathbf{b}_i$, as stated in Theorem 2(ii). Since $\rho_i$ and $\mathbf{b}_i$ represent loadings for all possible factors, this result is useful because these coefficients can be jointly estimated and inferred in many economic contexts, such as financial networks, as described in Wang and Shojaie (2021).

Next, we provide the joint limiting distributions of the shrinkage estimators. For $i = 1, ..., N$, define

$$\boldsymbol{\Sigma}_{f\varepsilon_i}(0,0) = \mathrm{Cov}(\mathbf{f}_t\varepsilon_{i,t}, \mathbf{f}_t\varepsilon_{i,t}), \ \ \boldsymbol{\Sigma}_{f\varepsilon_i}(1,0) = \mathrm{Cov}(\mathbf{f}_t\varepsilon_{i,t}, \mathbf{f}_{t-1}\varepsilon_{i,t}), \ \ \boldsymbol{\Omega}_{f\varepsilon_i}(0,0) = \mathrm{Cov}(\mathbf{f}_{t-1}\varepsilon_{i,t}, \mathbf{f}_{t-1}\varepsilon_{i,t}),$$

$$\boldsymbol{\Sigma}_{f\varepsilon_i}(k,j) = \mathrm{Cov}(\mathbf{f}_{t+j}\varepsilon_{i,t+j}, \mathbf{f}_{t-k}\varepsilon_{i,t}) + \mathrm{Cov}(\mathbf{f}_t\varepsilon_{i,t}, \mathbf{f}_{t-k+j}\varepsilon_{i,t+j}), \ j \geq 1, k \geq 0,$$

$$\boldsymbol{\Omega}_{f\varepsilon_i}(0,j) = \mathrm{Cov}(\mathbf{f}_{t-1+j}\varepsilon_{i,t+j}, \mathbf{f}_{t-1}\varepsilon_{i,t}) + \mathrm{Cov}(\mathbf{f}_{t-1}\varepsilon_{i,t}, \mathbf{f}_{t-1+j}\varepsilon_{i,t+j}), \quad j \geq 1,$$

$$\boldsymbol{\Sigma}_{f\varepsilon_i}(0) = \sum_{j=0}^{\infty}\boldsymbol{\Sigma}_{f\varepsilon_i}(0,j), \ \boldsymbol{\Sigma}_{f\varepsilon_i}(1) = \sum_{j=0}^{\infty}\boldsymbol{\Sigma}_{f\varepsilon_i}(1,j), \ \boldsymbol{\Omega}_{f\varepsilon_i}(0) = \sum_{j=0}^{\infty}\boldsymbol{\Omega}_{f\varepsilon_i}(0,j).$$

Let

$$\mathbf{V}_i = \begin{pmatrix} \mathbf{w}_i'\boldsymbol{\Sigma}_{yf}\boldsymbol{\Sigma}_{yf}'\mathbf{w}_i + \mathbf{w}_i'\boldsymbol{\Sigma}_{yf}(1)\boldsymbol{\Sigma}_{yf}'(1)\mathbf{w}_i & \mathbf{w}_i'\boldsymbol{\Sigma}_{yf}\boldsymbol{\Sigma}_f + \mathbf{w}_i'\boldsymbol{\Sigma}_{yf}(1)\boldsymbol{\Sigma}_f'(1) \\ \boldsymbol{\Sigma}_f\boldsymbol{\Sigma}_{yf}'\mathbf{w}_i + \boldsymbol{\Sigma}_f(1)\boldsymbol{\Sigma}_{yf}'(1)\mathbf{w}_i & \boldsymbol{\Sigma}_f^2 + \boldsymbol{\Sigma}_f(1)\boldsymbol{\Sigma}_f'(1) \end{pmatrix} \tag{23}$$

and

$$\mathbf{U}_i = \begin{pmatrix} \boldsymbol{\Sigma}_{f\varepsilon_i}(0) & \boldsymbol{\Sigma}_{f\varepsilon_i}(1) \\ \boldsymbol{\Sigma}_{f\varepsilon_i}'(1) & \boldsymbol{\Omega}_{f\varepsilon_i}(0) \end{pmatrix}. \tag{24}$$

The following theorem establishes the joint asymptotic normality of the estimators.

**Theorem 3.** *Let Assumptions 1 − 6 hold. If $N = o(T)$ and $\lambda_i \to 0$, we have*

$$\sqrt{T}\mathbf{V}_i(\widehat{\boldsymbol{\beta}}_i - \boldsymbol{\beta}_i) \longrightarrow_d N(\mathbf{0}, \mathbf{X}_i'\mathbf{U}_i\mathbf{X}_i),$$

*for $i = 1, ..., N$ as $N, T \to \infty$, where $\mathbf{U}_i$ and $\mathbf{V}_i$ are defined in (24) and (23), respectively.*

From Theorem 3, we see that the Yule-Walker estimators obtained in (12) are asymptotically normal when the dimension $N$ diverges. The convergence rate is the standard $\sqrt{T}$ under the

17

assumption that $N/T \to 0$, which is a similar requirement in spatial panel dynamic models; see Yu et al. (2008), among others. The condition $N/T \to 0$ is weaker than the one in Gao et al. (2019), where $N/\sqrt{T} \to 0$ is required, because we assume the dimension of $\mathbf{f}_t$ is $K$, a finite integer. The convergence of $\widehat{\boldsymbol{\Sigma}}_{yf}(k)$ to $\boldsymbol{\Sigma}_{yf}(k)$ only requires $N/T \to 0$, whereas the convergence of $\widehat{\boldsymbol{\Sigma}}_y(k) = \frac{1}{T} \sum_{t=k+1}^{T} \mathbf{y}_t \mathbf{y}'_{t-k}$ to its population version requires $N/\sqrt{T} \to 0$ as stated in Gao et al. (2019). By the form of $\mathbf{X}_i$ in Assumption 6, we can show that

$$\mathbf{X}'_i \mathbf{U}_i \mathbf{X}_i = \begin{pmatrix} \boldsymbol{\Sigma}_{i,11} & \boldsymbol{\Sigma}_{i,12} \\ \boldsymbol{\Sigma}_{i,21} & \boldsymbol{\Sigma}_{i,22} \end{pmatrix}, \tag{25}$$

where

$$\begin{aligned} \boldsymbol{\Sigma}_{i,11} =& \mathbf{w}'_i \boldsymbol{\Sigma}_{yf} \boldsymbol{\Sigma}_{f\varepsilon_i}(0) \boldsymbol{\Sigma}'_{yf} \mathbf{w}_i + \mathbf{w}'_i \boldsymbol{\Sigma}_{yf} \boldsymbol{\Sigma}_{f\varepsilon_i}(1) \boldsymbol{\Sigma}'_{yf}(1) \mathbf{w}_i + \mathbf{w}'_i \boldsymbol{\Sigma}_{yf}(1) \boldsymbol{\Sigma}'_{f\varepsilon_i}(1) \boldsymbol{\Sigma}'_{yf} \mathbf{w}_i \\ &+ \mathbf{w}'_i \boldsymbol{\Sigma}_{yf}(1) \boldsymbol{\Omega}_{f\varepsilon_i}(0) \boldsymbol{\Sigma}'_{yf}(1) \mathbf{w}_i, \end{aligned}$$

$$\begin{aligned} \boldsymbol{\Sigma}_{i,22} =& \boldsymbol{\Sigma}_f \boldsymbol{\Sigma}_{f\varepsilon_i}(0) \boldsymbol{\Sigma}_f + \boldsymbol{\Sigma}_f \boldsymbol{\Sigma}_{f\varepsilon_i}(1) \boldsymbol{\Sigma}'_f(1) + \boldsymbol{\Sigma}_f(1) \boldsymbol{\Sigma}'_{f\varepsilon_i}(1) \boldsymbol{\Sigma}_f + \boldsymbol{\Sigma}_f(1) \boldsymbol{\Omega}_{f\varepsilon_i}(0) \boldsymbol{\Sigma}'_f(1), \end{aligned}$$

$$\begin{aligned} \boldsymbol{\Sigma}_{i,12} =& \mathbf{w}'_i \boldsymbol{\Sigma}_{yf} \boldsymbol{\Sigma}_{f\varepsilon_i}(0) \boldsymbol{\Sigma}_f + \mathbf{w}'_i \boldsymbol{\Sigma}_{yf} \boldsymbol{\Sigma}_{f\varepsilon_i}(1) \boldsymbol{\Sigma}'_f(1) + \mathbf{w}'_i \boldsymbol{\Sigma}_{yf}(1) \boldsymbol{\Sigma}'_{f\varepsilon_i}(1) \boldsymbol{\Sigma}_f \\ &+ \mathbf{w}'_i \boldsymbol{\Sigma}_{yf}(1) \boldsymbol{\Omega}_{f\varepsilon_i}(0) \boldsymbol{\Sigma}'_f(1), \end{aligned}$$

and $\boldsymbol{\Sigma}_{i,21} = \boldsymbol{\Sigma}'_{i,12}$. These matries can all be estimated from the data.

Finally, we turn to the case when the factors are latent. We need to make two more assumptions to establish the uniform convergence and the limiting distributions of the estimated factors.

**Assumption 7.** *$\mathbf{f}_t$ and $\boldsymbol{\varepsilon}_t$ are sub-exponentially distributed in the sense that*

$$P(|\mathbf{v}'_1 \mathbf{f}_t| > x) \le C \exp(-Cx), \text{ and } P(|\mathbf{v}'_2 \boldsymbol{\varepsilon}_t| > x) \le C \exp(-Cx),$$

*for any $x > 0$, where $\|\mathbf{v}_1\|_2 = 1$ and $\|\mathbf{v}_2\|_2 = 1$ are any two constant vectors.*

**Assumption 8.** *For each $t = 1, ..., T$, as $N \to \infty$,*

$$\frac{1}{\sqrt{N}} \sum_{i=1}^{N} \mathbf{p}_i \varepsilon_{i,t} \longrightarrow_d N(0, \boldsymbol{\Gamma}_t),$$

*where $\mathbf{p}_i$ is the ith column of $\boldsymbol{\Lambda}'(\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1}$, and $\boldsymbol{\Gamma}_t = \lim_{N\to\infty} \frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{N} \mathbf{p}_i \mathbf{p}'_j E(\varepsilon_{i,t} \varepsilon_{j,t})$ in probability.*

Assumption 7 is commonly used in the statistical and econometrics literature to establish uniform convergence. The sub-exponential distribution is a broader class of distributions than the sub-Gaussian distribution and includes the uniform distribution over every convex body, following

18

the Brunn-Minkowski inequality. For further details, see, for example, Vershynin (2018). Assumption 8 is similar to Assumption F(3) in Bai (2003), which is used to establish the limiting distribution of the estimated factors.

We first state the convergence of the estimated loading matrix below, where we introduce a rotational matrix $\mathbf{H}_{NT}$ in the proof of the following theorem. This approach differs from the techniques used in Lam et al. (2011).

**Theorem 4.** *Let Assumptions $1 - 6$ hold. If $N = o(T)$, then there exists an invertible matrix $\mathbf{H}_{NT}$ such that*

$$\frac{1}{\sqrt{N}}\|\widehat{\boldsymbol{\Lambda}} - \boldsymbol{\Lambda}\mathbf{H}'_{NT}\|_F = O_p(\frac{1}{\sqrt{T}}).$$

**Remark 1.** *(i) Unlike the proof in Lam et al. (2011) where a matrix perturbation theory is used to show the convergence of the estimated loading matrix, we developed a new approach in the Appendix to show the convergence rate of $\widehat{\boldsymbol{\Lambda}}$. One of the advantages of the new approach is that we can specify the rotational matrix $\mathbf{H}_{NT}$ which is defined as*

$$\mathbf{H}'_{NT} = \sum_{k=1}^{k_0}\mathbf{G}_{1,k}\mathbf{G}'_{1,k}\boldsymbol{\Lambda}'\widehat{\boldsymbol{\Lambda}}\widehat{\mathbf{V}}_{NT}^{-1}, \text{ where } \mathbf{G}_{1,k} = \frac{1}{T}\sum_{t=k+1}^{T}(\mathbf{f}_t\mathbf{f}'_{t-k}\boldsymbol{\Lambda}' + \mathbf{f}_t\boldsymbol{\xi}'_{t-k}),$$

*and $\widehat{\mathbf{V}}_{NT} \in R^r$ is a diagonal matrix with diagonal elements being the top $K$ eigenvalues of $\widehat{\mathbf{M}}$. See the proof of Theorem 4 in the online Appendix for details.*
*(ii) Note that we impose that $\boldsymbol{\Lambda}'\boldsymbol{\Lambda}/N = \mathbf{I}_r$, whereas Lam et al. (2011) assumes that $\boldsymbol{\Lambda}'\boldsymbol{\Lambda} = \mathbf{I}_r$. Therefore, the convergence rate is the same as the one in Theorem 1 of Lam et al. (2011), where we assume $\delta = 0$ in our paper, corresponding to the case of strong factors.*

Next, we establish the uniform convergence of the estimated factors and the corresponding limiting distributions.

**Theorem 5.** *Let Assumptions $1 - 6$ hold.*
*(i) If $\boldsymbol{\varepsilon}_t$ and $\mathbf{f}_t$ are sub-exponentially distributed as in Assumption 7, then there exists an invertible matrix $\mathbf{K}_{NT} \in \mathbb{R}^r$ such that*

$$\max_{1 \leq t \leq T}\|\widehat{\mathbf{f}}_t - \mathbf{K}_{NT}\mathbf{f}_t\|_2 = O_p\{(\frac{1}{\sqrt{N}} + \frac{1}{\sqrt{T}})\log(T)\}.$$

*(ii) Let Assumption 8 also hold. If $N = o(T)$, then there exists an invertible matrix $\mathbf{K}_{NT} \in \mathbb{R}^r$ and its limit $\mathbf{H} \in \mathbb{R}^r$ such that*

$$\sqrt{N}(\widehat{\mathbf{f}}_t - \mathbf{K}_{NT}\mathbf{f}_t) = \mathbf{H}\frac{1}{\sqrt{N}}\sum_{i=1}^{N}\mathbf{p}_i\varepsilon_{i,t} + o_p(1) \longrightarrow_d N(0, \mathbf{H}\boldsymbol{\Gamma}_t\mathbf{H}'),$$

*where $\mathbf{p}_i$ is the ith column of $\boldsymbol{\Lambda}'(\mathbf{I}_N - \mathbf{D}(\boldsymbol{\rho})\mathbf{W})^{-1}$, $\mathbf{H}$ is the limit of $\mathbf{H}_{NT}$ as shown in Lemma 2 of the online Appendix, and $\boldsymbol{\Gamma}_t$ is defined as in Assumption 8.*

**Remark 2.** *(I) A remarkable feature in Theorem 5 is that we only require $N/T \to 0$, and the asymptotic normality of $\mathbf{f}_t$ can still be achieved.*

*(ii) Note that we adopt the matrix $\mathbf{K}_{NT}$ as a rotational matrix for $\mathbf{f}_t$, which is defined as*

$$\mathbf{K}_{NT} = \frac{1}{N}\widehat{\boldsymbol{\Lambda}}'\boldsymbol{\Lambda}.$$

*See the proof of Theorem 5 in the Appendix. In fact, according to Lemma 1 of the Appendix, we may replace $\mathbf{K}_{NT}$ by $\mathbf{H}_{NT}$, and the results in Theorem 5 still hold. This can be shown by rewriting the term $\frac{1}{N}\widehat{\boldsymbol{\Lambda}}'\boldsymbol{\Lambda}\mathbf{f}_t$ in (IA.16) as*

$$\frac{1}{N}(\widehat{\boldsymbol{\Lambda}} - \boldsymbol{\Lambda}\mathbf{H}'_{NT})'\boldsymbol{\Lambda}\mathbf{f}_t + \mathbf{H}_{NT}\mathbf{f}_t,$$

*where the first term is still asymptotically negligible. However, we do not adopt this formula since it will introduce a bias term in establishing the limiting distributions of the $\widehat{\mathbf{f}}_t$. Nevertheless, it is not hard to show that $\mathbf{K}_{NT}$ and $\mathbf{H}_{NT}$ have the same limit as $N, T \to \infty$.*

Furthermore, we study the limiting distributions of the estimated parameters in (19). Similar to the case when the factors are observable, we provide some notation used in the following Theorem. Let

$$\mathbf{V}_i^H = \begin{pmatrix} \mathbf{w}_i'\boldsymbol{\Sigma}_{yf}\boldsymbol{\Sigma}'_{yf}\mathbf{w}_i + \mathbf{w}_i'\boldsymbol{\Sigma}_{yf}(1)\boldsymbol{\Sigma}'_{yf}(1)\mathbf{w}_i & \mathbf{w}_i'\boldsymbol{\Sigma}_{yf}\boldsymbol{\Sigma}_f\mathbf{H}' + \mathbf{w}_i'\boldsymbol{\Sigma}_{yf}(1)\boldsymbol{\Sigma}'_f(1)\mathbf{H}' \\ \mathbf{H}\boldsymbol{\Sigma}_f\boldsymbol{\Sigma}'_{yf}\mathbf{w}_i + \mathbf{H}\boldsymbol{\Sigma}_f(1)\boldsymbol{\Sigma}'_{yf}(1)\mathbf{w}_i & \mathbf{H}\boldsymbol{\Sigma}_f^2\mathbf{H}' + \mathbf{H}\boldsymbol{\Sigma}_f(1)\boldsymbol{\Sigma}'_f(1)\mathbf{H}' \end{pmatrix} \tag{26}$$

and

$$\mathbf{U}_i^H = \begin{pmatrix} \mathbf{H}\boldsymbol{\Sigma}_{f\varepsilon_i}(0)\mathbf{H}' & \mathbf{H}\boldsymbol{\Sigma}_{f\varepsilon_i}(1)\mathbf{H}' \\ \mathbf{H}\boldsymbol{\Sigma}'_{f\varepsilon_i}(1)\mathbf{H}' & \mathbf{H}\boldsymbol{\Omega}_{f\varepsilon_i}(0)\mathbf{H}' \end{pmatrix}. \tag{27}$$

The following theorem establishes the asymptotic normality of the estimators in (19) when the factors are latent and the dimension $N$ is diverging.

**Theorem 6.** *Let Assumptions $1 - 8$ hold.*
*(i) If $N = o(T)$ and $\sqrt{T} = o(N)$, then there exists an invertible matrix $\mathbf{K}_{NT} \in R^r$ such that*

$$\widetilde{\boldsymbol{\beta}}_i(\lambda_i) - \widetilde{\mathbf{X}}_i(\lambda_i)^{-1}\widetilde{\mathbf{X}}'_i\widetilde{\mathbf{X}}_i\mathbf{K}^*_{NT}\boldsymbol{\beta}_i = O_p(T^{-1/2}),$$

*where $\widetilde{\mathbf{X}}_i(\lambda_i) = \widetilde{\mathbf{X}}'_i\widetilde{\mathbf{X}}_i + \lambda_i\mathbf{I}_{K+1}$.*
*(ii) If $N = o(T)$ and $\sqrt{T} = o(N)$, let $\lambda_i \to 0$, there exists an invertible matrix $\mathbf{K}_{NT} \in R^r$ such that*

$$\sqrt{T}\mathbf{V}_i^H(\widetilde{\boldsymbol{\beta}}_i - \mathbf{K}^*_{NT}\boldsymbol{\beta}_i) \longrightarrow_d N(\mathbf{0}, \mathbf{X}_i^{H'}\mathbf{U}_i^H\mathbf{X}_i^H),$$

*for $i = 1, ..., N$ as $T \to \infty$, where $\mathbf{K}^*_{NT} = diag(1, (\mathbf{K}'_{NT})^{-1})$ is a block-diagonal matrix, and $\mathbf{U}_i^H$*

and $\mathbf{V}_i^H$ are defined in (27) and (26), respectively, and

$$\mathbf{X}_i^H = \begin{pmatrix} \mathbf{H}\mathbf{\Sigma}'_{yf}\mathbf{w}_i & \mathbf{H}\mathbf{\Sigma}_f\mathbf{H}' \\ \mathbf{H}\mathbf{\Sigma}'_{yf}(1)\mathbf{w}_i & \mathbf{H}\mathbf{\Sigma}'_f(1)\mathbf{H}' \end{pmatrix}.$$

**Remark 3.** *(i) From Theorem 6, we see that the convergence rate is still the standard $\sqrt{T}$, which is the same as that in Theorem 3 when the factors are observable. On the other hand, we note that the scalar coefficient can be uniquely determined, but the coefficient vector $\mathbf{b}_i$ can be estimated up to a rotational matrix $\mathbf{K}_{NT}$, which is reasonable due to the identification issue in the factor analysis. (ii) Recall that this is a two-step procedure. The statistical inference is usually difficult to establish in the second step because the errors incurred in the first step sometimes create a biased term. As discussed in Remark 2(ii), we adopt a rotational matrix $\mathbf{K}_{NT}$ instead of $\mathbf{H}_{NT}$ in Theorems 5 and 6 such that the bias term can be erased, although $\mathbf{K}_{NT}$ and $\mathbf{H}_{NT}$ have the same limit. See the proof of Theorem 6 in the online Appendix for details.*

It can be easily shown that the variance term in Theorem 6(ii) can be expressed as

$$\mathbf{X}_i^{H\prime}\mathbf{U}_i^H\mathbf{X}_i^H = \begin{pmatrix} \mathbf{\Sigma}_{i,11}^H & \mathbf{\Sigma}_{i,12}^H \\ \mathbf{\Sigma}_{i,21}^H & \mathbf{\Sigma}_{i,22}^H \end{pmatrix},$$

where

$$\begin{aligned}
\mathbf{\Sigma}_{i,11}^H =\,& \mathbf{w}'_i\mathbf{\Sigma}_{yf}\mathbf{\Sigma}_{f\varepsilon_i}(0)\mathbf{\Sigma}'_{yf}\mathbf{w}_i + \mathbf{w}'_i\mathbf{\Sigma}_{yf}\mathbf{\Sigma}_{f\varepsilon_i}(1)\mathbf{\Sigma}'_{yf}(1)\mathbf{w}_i + \mathbf{w}'_i\mathbf{\Sigma}_{yf}(1)\mathbf{\Sigma}'_{f\varepsilon_i}(1)\mathbf{\Sigma}'_{yf}\mathbf{w}_i \\
& + \mathbf{w}'_i\mathbf{\Sigma}_{yf}(1)\mathbf{\Omega}_{f\varepsilon_i}(0)\mathbf{\Sigma}'_{yf}(1)\mathbf{w}_i,
\end{aligned}$$

$$\mathbf{\Sigma}_{i,22}^H = \mathbf{H}\mathbf{\Sigma}_f\mathbf{\Sigma}_{f\varepsilon_i}(0)\mathbf{\Sigma}_f\mathbf{H}' + \mathbf{H}\mathbf{\Sigma}_f\mathbf{\Sigma}_{f\varepsilon_i}(1)\mathbf{\Sigma}'_f(1)\mathbf{H}' + \mathbf{H}\mathbf{\Sigma}_f(1)\mathbf{\Sigma}'_{f\varepsilon_i}(1)\mathbf{\Sigma}_f\mathbf{H}' + \mathbf{H}\mathbf{\Sigma}_f(1)\mathbf{\Omega}_{f\varepsilon_i}(0)\mathbf{\Sigma}'_f(1)\mathbf{H}',$$

$$\begin{aligned}
\mathbf{\Sigma}_{i,12}^H =\,& \mathbf{w}'_i\mathbf{\Sigma}_{yf}\mathbf{\Sigma}_{f\varepsilon_i}(0)\mathbf{\Sigma}_f\mathbf{H}' + \mathbf{w}'_i\mathbf{\Sigma}_{yf}\mathbf{\Sigma}_{f\varepsilon_i}(1)\mathbf{\Sigma}'_f(1)\mathbf{H}' + \mathbf{w}'_i\mathbf{\Sigma}_{yf}(1)\mathbf{\Sigma}'_{f\varepsilon_i}(1)\mathbf{\Sigma}_f\mathbf{H}' \\
& + \mathbf{w}'_i\mathbf{\Sigma}_{yf}(1)\mathbf{\Omega}_{f\varepsilon_i}(0)\mathbf{\Sigma}'_f(1)\mathbf{H}',
\end{aligned}$$

and $\mathbf{\Sigma}_{i,21}^H = \mathbf{\Sigma}_{i,12}^H{}'$.

The consistency of the estimated number of factors using the information criterion in (20) or the ratio-based method in (21) can be established by a standard argument as that in Bai and Ng (2002) or Ahn and Horenstein (2013). We omit the details.

**Remark 4.** *In the estimation procedure above, we primarily focus on the augmented method by stacking factor lags for $k = 0$ and $k = 1$. In fact, ridge regression can be applied by taking any finite number of lagged factors in the Yule-Walker estimation. The theory can be established in a similar way.*

**Remark 5.** *The QMLE method proposed in Aquaro et al. (2021) and Hu et al. (2023) can yield pointwise consistent estimators but is feasible only when the dimension $N$ is small and fixed. Additionally, they only focus on cases when the factors are observable. As $N$ increases, additional bias can arise, and the asymptotic results do not hold anymore (see Remarks 6–7 in Aquaro et al. (2021)). Moreover, the computational cost of the QMLE method becomes prohibitive for large $N$. In contrast, the proposed generalized Yule-Walker method is designed to handle scenarios with large or diverging $N$ while remaining computationally efficient. Simulations and real data analyses in Sections 5-6 show that the proposed method can even outperform the QML approach, achieving smaller out-of-sample forecasting errors.*

## 5 Simulation Studies

In this section, we use Monte Carlo simulations to evaluate the performance of the proposed methodology across a spectrum of finite samples.

Consider the model in Section 3 with common factors generated from a VAR(1) process $\mathbf{f}_t = \mathbf{\Phi f}_{t-1} + \boldsymbol{\eta}_t$. Here, $\mathbf{\Phi}$ is a diagonal matrix, with entries independently sampled from a uniform distribution $U(0.5, 0.9)$ and the error term $\boldsymbol{\eta_t} \sim N(0, \mathbf{I}_K)$. For each realization of $\mathbf{y}_t$, the elements of the loading matrix $\mathbf{B}$ are independently drawn from $U(-2, 2)$, and the idiosyncratic error term $\boldsymbol{\varepsilon}_t$ is generated from $N(0, I_N)$. The spatial $\boldsymbol{\rho}$ is sampled independently from a power-law distribution with an exponent $\alpha = 5$. To construct the spatial matrices, the $q$ neighboring off-diagonal elements are set to 1 and the diagonal elements are 0, followed by row normalization to ensure each row sums to 1. We set $q = 3$ and the true number of factors $K = 3$, with dimension $N = 25, 50, 100, 200$, and sample size $T = 50, 100, 200, 400, 1000$. We use 1000 replications for each configuration of $(T, N)$. To make the results below replicable, the seed is set to be `1234` in the `R` programming.

We first examine the joint convergence properties of $\widehat{\boldsymbol{\beta}}_i$ established in Theorem 2. To evaluate its overall estimation accuracy, we use the root-mean-square error (RMSE), defined as

$$
\text{RMSE}_{\widehat{\beta}} = \begin{cases} \left(\frac{1}{N}\sum_{i=1}^{N}\|(\widehat{\mathbf{X}}_i'\widehat{\mathbf{X}}_i)(\widehat{\boldsymbol{\beta}}_i - \boldsymbol{\beta}_i)\|_2^2\right)^{1/2}, & \text{if } \lambda_i \to 0 \\ \left(\frac{1}{N}\sum_{i=1}^{N}\|\widehat{\boldsymbol{\beta}}_i(\lambda_i) - \widehat{\mathbf{X}}_i(\lambda_i)^{-1}\widehat{\mathbf{X}}_i\widehat{\mathbf{X}}_i'\boldsymbol{\beta}_i\|_2^2\right)^{1/2}, & \text{otherwise.} \end{cases} \tag{28}
$$

Here, $\lambda_i$ is the ridge penalty parameter applied to the Yule-Walker equations for each sample. We examine two cases: a relatively large $\lambda_i = 10^{-3}$ and a much smaller $\lambda_i = 10^{-9}$. When $\lambda_i \to 0$ (e.g., $\lambda_i = 10^{-9}$), the estimator closely resembles that of the ordinary least squares (OLS) estimation, but we set $\lambda_i = 10^{-9}$ to avoid singularity of $(\widehat{\mathbf{X}}_i'\widehat{\mathbf{X}}_i)$. Figure 2(a) and (b) present the boxplots of the RMSEs of $\widehat{\boldsymbol{\beta}}(\lambda_i)$'s (denoted by $\text{RMSE}_{\widehat{\beta}}$) and $\widehat{\rho}_i(\lambda_i)$'s (denoted by $\text{RMSE}_{\widehat{\rho}}$), respectively, computed using the second formula in (28). From Figure 2, we see that the $\text{RMSE}_{\widehat{\beta}}$ and $\text{RMSE}_{\widehat{\rho}}$ decrease as the sample size $T$ increases, which is in agreement with the theoretical results in Theorem 2. Similar patterns can also be found in Figure 3 for $\lambda \to 0$ using the RMSE defined in the first line of (28).
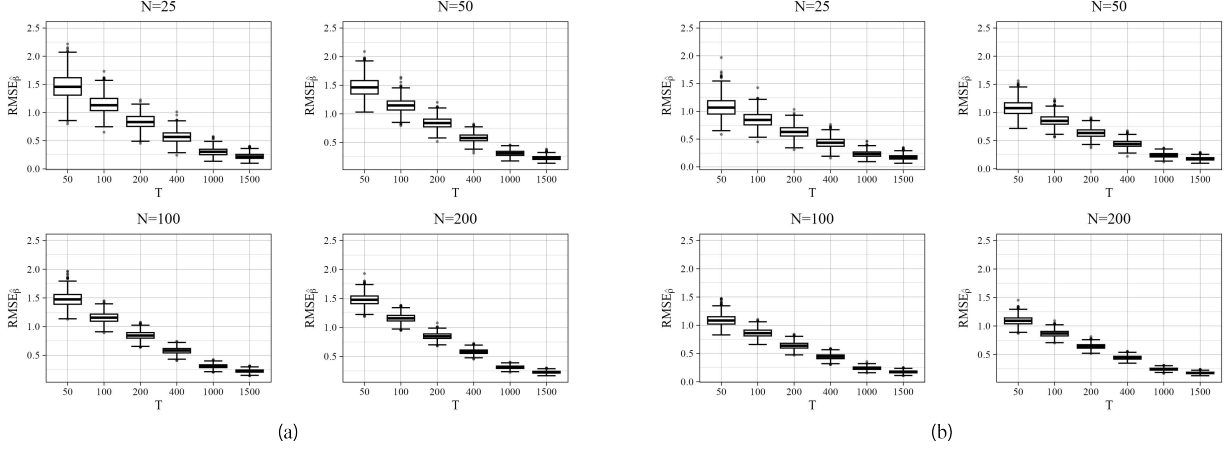
**Figure 2:** Boxplots of estimator convergence for Eq (11) with a fixed large ridge penalty parameter $(\lambda_i = 10^{-3})$, where $N$ and $T$ denote the dimension and sample size, respectively. (a) shows the joint estimation performance of $\widehat{\boldsymbol{\beta}}_i$ measured by $\mathrm{RMSE}_{\widehat{\beta}}$, and (b) shows the estimation performance of the spatial parameter $\widehat{\boldsymbol{\rho}}$ measured by $\mathrm{RMSE}_{\widehat{\rho}}$.

**Table 1:** Coverage rates for $\mathbf{V}_i(\widehat{\boldsymbol{\beta}}_i - \boldsymbol{\beta}_i)$ across different significance levels.

| Significance | Coverage | | | | | |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| **0.1** | 0.801(0.013) | 0.852(0.011) | 0.941(0.007) | 0.966(0.006) | 0.989(0.003) | 0.998(0.001) |
| **0.05** | 0.842(0.012) | 0.891(0.010) | 0.957(0.006) | 0.983(0.004) | 0.997(0.002) | 1.000(0.000) |
| **0.01** | 0.898(0.010) | 0.935(0.008) | 0.980(0.004) | 0.993(0.003) | 0.998(0.001) | 1.000(0.000) |
| $T$ | **250** | **500** | **1000** | **2000** | **3000** | **5000** |

*Note:* This table shows the coverage rate of the first component of $\mathbf{V}_1(\widehat{\boldsymbol{\beta}}_1 - \boldsymbol{\beta}_1)$ within the confidence intervals from the theoretical distribution in Theorem 3. The theoretical distribution has a zero mean and a variance equal to the first diagonal element of $\mathbf{X}_1'\mathbf{U}_1\mathbf{X}_1$. $T$ is the sample size. The results are based on 1000 iterations with a cross-sectional dimension $N = 25$.

To assess the distributional properties of the estimates in Theorem 3, we present histograms of the first component of $\widehat{\mathbf{V}}_i(\widehat{\boldsymbol{\beta}}_i - \boldsymbol{\beta}_i)$, for $i = 1, 3, 5, 7$ and 9 in Figure 4, along with their theoretical density curves computed from the limiting distribution in Theorem 3, where $\widehat{\mathbf{V}}_i$ is the sample estimator for $\mathbf{V}_i$ defined in Eq (23). The histograms and their corresponding QQ-plots in Figure 4 suggest that the entries of $\widehat{\mathbf{V}}_i(\widehat{\boldsymbol{\beta}}_i - \boldsymbol{\beta}_i)$ asymptotically follow a normal distribution, which aligns with our theoretical results. Furthermore, in Table 1, we evaluate the asymptotic properties by reporting the coverage rates of the estimators in Eq (11) under varying significance levels and sample sizes $T$. As $T$ increases, the coverage rates exhibit a clear improvement, consistent with the theoretical results in Theorem 3.

As a shrinkage-based approach, it is interesting to directly assess the estimation accuracy of $\boldsymbol{\beta}_i$.
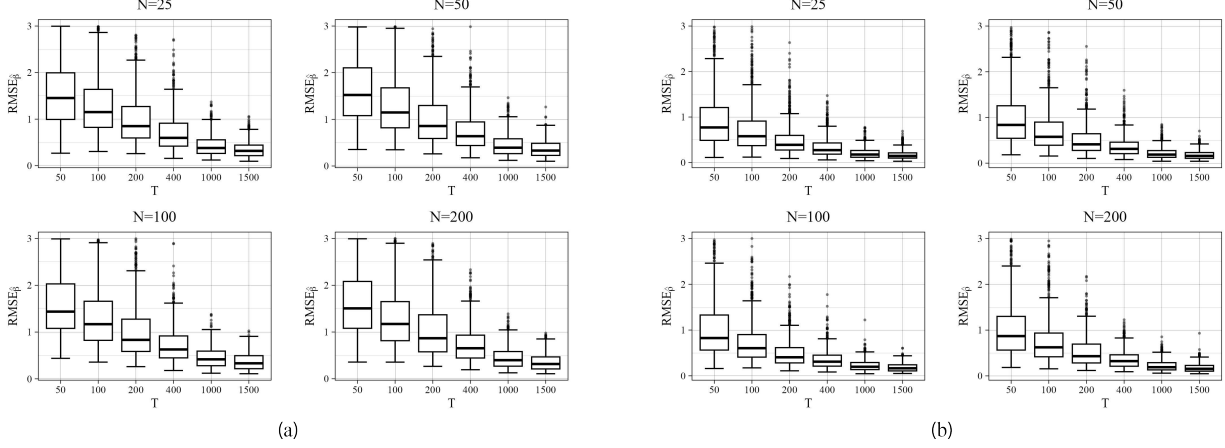
**Figure 3:** Boxplots of the joint convergence error for model (11) with $\lambda = 0$. The statistics are defined in Figure 2, based on the first line of (28).
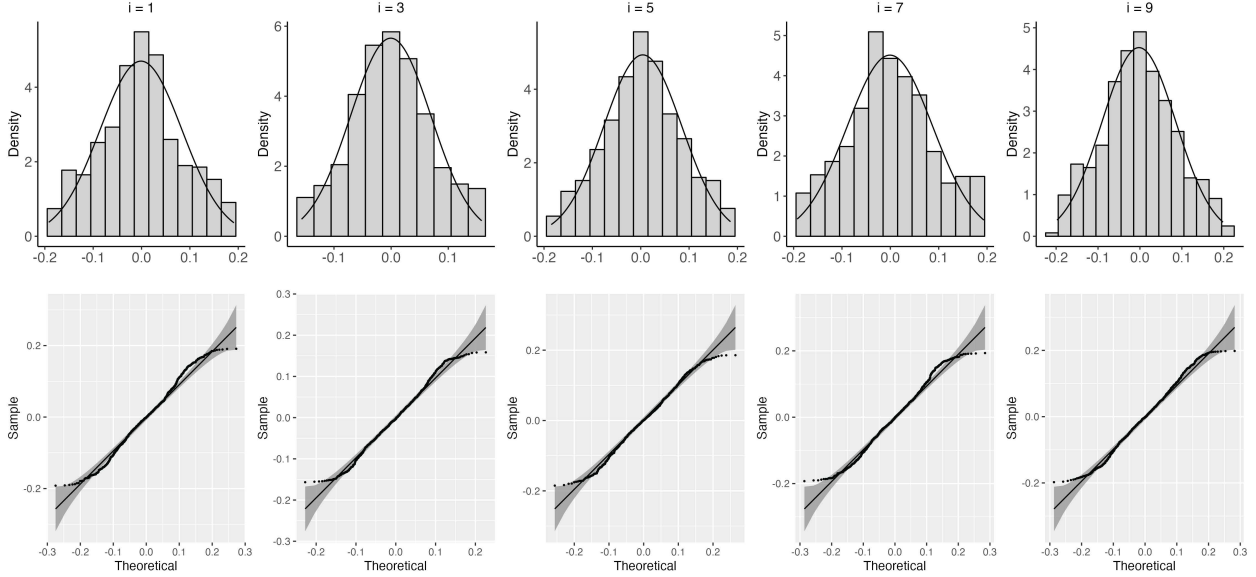
.



**Figure 4:** Histograms of the five spatial coefficient estimates and their corresponding empirical and theoretical distribution plots for Eq (11). The histograms show the distribution of the first component of $\mathbf{V}_i(\widehat{\boldsymbol{\beta}}_i - \boldsymbol{\beta}_i)$ over 1000 iterations, under the setting of $N = 25$ and $T = 1500$. The superimposed normal curve represents the theoretical distribution from Theorem 3, with mean 0 and variance given by the first diagonal element of $\mathbf{X}_i'\mathbf{U}_i\mathbf{X}_i$. Here, $i = 1, 3, 5, 7, 9$ correspond to the 1st, 3rd, 5th, 7th, and 9th component in a dataset of dimension $N = 25$.

We define the coefficient error (CE) as

$$\text{CE}_{\widehat{\beta}} = \left( \frac{1}{N} \sum_{i=1}^{N} \|(\widehat{\boldsymbol{\beta}}_i(\lambda_i) - \boldsymbol{\beta}_i)\|_2^2 \right)^{1/2}, \tag{29}$$

which measures the deviation of $\widehat{\boldsymbol{\beta}}_i$ from the true parameter $\boldsymbol{\beta}_i$, thereby capturing the overall

24

estimation error. Table 2 reports the coefficient error for the ridge regression estimates $\widehat{\boldsymbol{\beta}}_i(\lambda_i)$ when factors are observed. The results indicate that ridge regression ($\lambda_i = 10^{-3}$) yields lower error and variance compared to OLS ($\lambda_i = 10^{-9}$). Additionally, the estimation errors are reduced when stacking the cases of $k = 0$ and $k = 1$ together in the Yule-Walker equations, compared to relying solely on $k = 0$.

**Table 2:** Comparison of CE for ridge regression estimators with observed factors across different penalized parameter $\lambda_i$ and the lagging factor impact.

| N | T | $k = 1$ $(\lambda_i = 10^{-9})$ | | $k = 1$ $(\lambda_i = 10^{-3})$ | | $k = 0$ $(\lambda_i = 10^{-3})$ | |
|---|---|---|---|---|---|---|---|
| | | $\text{CE}_{\widehat{\beta}}$ | $\text{CE}_{\widehat{\rho}}$ | $\text{CE}_{\widehat{\beta}}$ | $\text{CE}_{\widehat{\rho}}$ | $\text{CE}_{\widehat{\beta}}$ | $\text{CE}_{\widehat{\rho}}$ |
| | 50 | 1.394(0.495) | 0.291(0.276) | 1.021(0.263) | 0.135(0.102) | 1.251(0.122) | 0.165(0.132) |
| | 100 | 1.405(0.510) | 0.357(0.326) | 1.012(0.260) | 0.124(0.101) | 1.238(0.129) | 0.164(0.136) |
| 25 | 200 | 1.444(0.509) | 0.488(0.379) | 0.971(0.235) | 0.120(0.097) | 1.239(0.125) | 0.167(0.131) |
| | 400 | 1.562(0.533) | 0.662(0.399) | 0.936(0.220) | 0.117(0.096) | 1.233(0.124) | 0.171(0.137) |
| | 1000 | 1.810(0.522) | 0.925(0.358) | 0.976(0.219) | 0.133(0.113) | 1.236(0.122) | 0.175(0.137) |
| | 1500 | 1.906(0.522) | 0.990(0.347) | 1.000(0.210) | 0.129(0.104) | 1.231(0.125) | 0.171(0.136) |
| | 50 | 1.424(0.444) | 0.328(0.248) | 1.009(0.210) | 0.154(0.089) | 1.219(0.086) | 0.203(0.124) |
| | 100 | 1.429(0.445) | 0.412(0.307) | 0.992(0.206) | 0.144(0.082) | 1.218(0.085) | 0.200(0.124) |
| 50 | 200 | 1.444(0.457) | 0.525(0.333) | 0.934(0.195) | 0.131(0.073) | 1.213(0.089) | 0.202(0.124) |
| | 400 | 1.587(0.479) | 0.713(0.349) | 0.900(0.182) | 0.130(0.079) | 1.215(0.084) | 0.202(0.124) |
| | 1000 | 1.838(0.460) | 0.947(0.310) | 0.931(0.174) | 0.133(0.079) | 1.205(0.087) | 0.192(0.116) |
| | 1500 | 1.918(0.445) | 1.010(0.290) | 0.966(0.170) | 0.144(0.089) | 1.208(0.086) | 0.198(0.125) |
| | 50 | 1.467(0.415) | 0.318(0.239) | 1.010(0.181) | 0.125(0.053) | 1.177(0.057) | 0.172(0.077) |
| | 100 | 1.471(0.400) | 0.407(0.271) | 0.995(0.173) | 0.120(0.048) | 1.172(0.059) | 0.168(0.075) |
| 100 | 200 | 1.502(0.414) | 0.556(0.325) | 0.916(0.158) | 0.104(0.043) | 1.175(0.061) | 0.163(0.074) |
| | 400 | 1.591(0.445) | 0.725(0.324) | 0.847(0.152) | 0.101(0.041) | 1.172(0.060) | 0.169(0.075) |
| | 1000 | 1.829(0.425) | 0.937(0.284) | 0.877(0.148) | 0.106(0.053) | 1.174(0.059) | 0.166(0.080) |
| | 1500 | 1.920(0.386) | 1.008(0.245) | 0.913(0.143) | 0.113(0.049) | 1.171(0.060) | 0.171(0.079) |
| | 50 | 1.443(0.431) | 0.273(0.238) | 0.988(0.195) | 0.086(0.042) | 1.174(0.038) | 0.120(0.064) |
| | 100 | 1.434(0.430) | 0.352(0.283) | 0.969(0.184) | 0.080(0.040) | 1.169(0.041) | 0.114(0.062) |
| 200 | 200 | 1.458(0.436) | 0.505(0.322) | 0.897(0.171) | 0.072(0.036) | 1.167(0.041) | 0.114(0.061) |
| | 400 | 1.566(0.462) | 0.680(0.346) | 0.843(0.150) | 0.068(0.034) | 1.165(0.041) | 0.113(0.061) |
| | 1000 | 1.834(0.453) | 0.938(0.303) | 0.868(0.145) | 0.072(0.037) | 1.165(0.042) | 0.116(0.062) |
| | 1500 | 1.910(0.427) | 0.995(0.274) | 0.909(0.140) | 0.076(0.039) | 1.167(0.041) | 0.117(0.062) |

*Note:* Here, $k = 1$ represents the application of the stacking strategy, which incorporates lagged factors $\mathbf{f}_{t-1}$ as instrumental variables, whereas $k = 0$ indicates the stacking strategy is not applied, relying exclusively on the contemporaneous factors $\mathbf{f}_t$, as detailed in Section 3.2.

Next, we investigate the performance of our proposed method in scenarios where the common factors are unobserved, focusing on the recovery of latent factors $\widehat{\mathbf{f}}_t$ and the estimation accuracy of the parameters $\boldsymbol{\rho}$ and $\mathbf{B}$ in Eq (14). Figure 5 illustrates the convergence behavior of the estimated loading matrix and latent factors using our method. In Panel (a) of Figure 5, it is evident that the loading matrix converges steadily as $T$ increases. On the other hand, for each fixed $N$, Panel (b) reveals that $\max_{1 \leq t \leq T} \|\widehat{\mathbf{f}}_t - \mathbf{K}_{NT}\mathbf{f}_t\|_2$ increases with $T$, which is reasonable since the uniform distance is measured over the entire $T$-period. However, for each fixed $T$, the uniform distance will become smaller as $N$ increases, which is in agreement with our theoretical results.

Now, we examine the asymptotic normality of the estimated latent factors. Figure 6 presents the histograms and QQ-plots of the first element of $\widehat{\mathbf{f}}_t - \mathbf{K}_{NT}\mathbf{f}_t$ for $t = 1, 3, 5, 7$ and 9 when $T = 1500$, which clearly show an asymptotic normality pattern across all settings. In addition, Table 3 presents the average coverage rates of the first component of $\widehat{\mathbf{f}}_t - \mathbf{K}_{NT}\mathbf{f}_t$ for $t \in \{1, 6, 11, ..., 96\}$, with a total of 20 factors, across different significance levels and values of $N$. As $N$ increases, the average coverage rates gradually improve, accompanied by reduced variance. These findings align with Theorem 4 and Theorem 5.
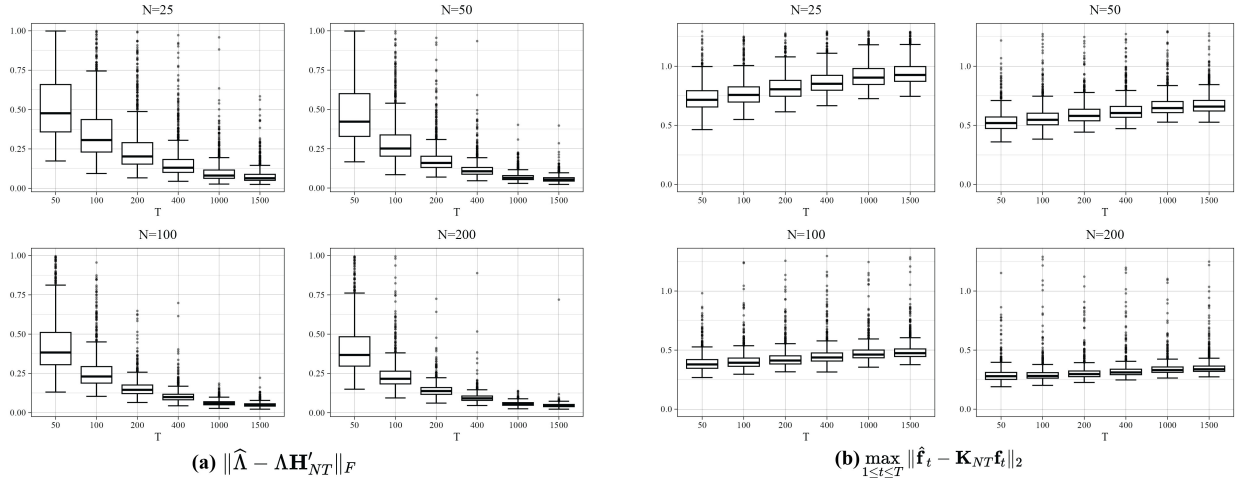


**(a)** $\|\widehat{\Lambda} - \Lambda\mathbf{H}'_{NT}\|_F$    **(b)** $\max_{1 \leq t \leq T} \|\widehat{\mathbf{f}}_t - \mathbf{K}_{NT}\mathbf{f}_t\|_2$

**Figure 5:** Boxplots of the convergence performance for the estimated loading matrices and latent factors in Eq (14).

With the estimated factors $\widehat{\mathbf{f}}_t$, Figure 7 demonstrates the boxplots of the RMSE of $\widetilde{\boldsymbol{\beta}}_i$ under increasing $T$, where the RMSE is similarly define as (28). The patterns of the RMSEs in Figure 7 are similar to those in Figure 3, and we omit the details here. To validate the distributional properties, Figure 8 displays the histogram and QQ-plot of the first component of $\mathbf{V}_i^H(\widetilde{\boldsymbol{\beta}}_i - \mathbf{K}_{NT}^*\boldsymbol{\beta}_i)$, for $i = 1, 3, 5, 7$ and 9. From Figure 7, we can see clearly an asymptotic normality pattern across all settings, which is in line with our theoretical results in Theorem 6. Moreover, we verify the coverage probabilities of the first component of $\mathbf{V}_1^H(\widetilde{\boldsymbol{\beta}}_1 - \mathbf{K}_{NT}^*\boldsymbol{\beta}_1)$ in Table 8, which are also in agreement with our theory.

Next, we present the coefficient error (CE) results for ridge estimators in the case of unknown factors in Table 5. From Table 5 we see that integrating the stacking strategy with proper ridge
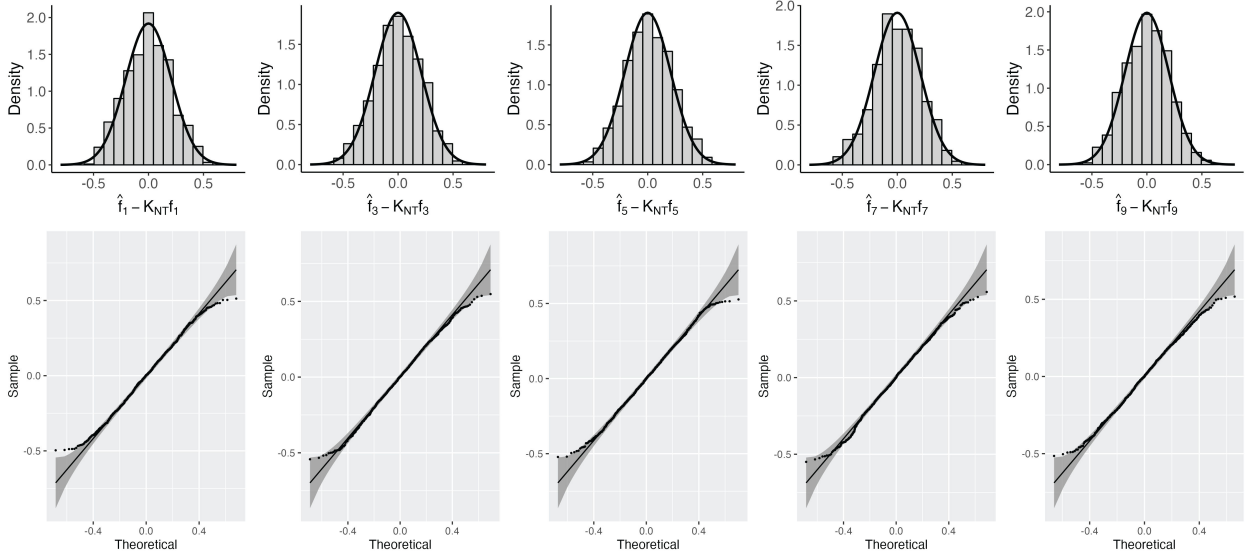
**Figure 6:** Histograms of $\widehat{\mathbf{f}}_t - \mathbf{K}_{NT}\mathbf{f}_t$ and their corresponding empirical and theoretical distribution plots for Model (14). The results are based on 1000 iterations, focusing on the first element of $\widehat{\mathbf{f}}_t - \mathbf{K}_{NT}\mathbf{f}_t$, for t = 1, 3, 5, 7, 9. The superimposed normal curves represent the theoretical distribution derived in Theorem 5(ii), with mean 0 and variance given by the first diagonal element of $\mathbf{H}\mathbf{\Gamma}_t\mathbf{H}'$. The simulation results are obtained under $(N, T) = (25, 1500)$.

**Table 3:** Coverage rates of $\widehat{\mathbf{f}}_t - \mathbf{K}_{NT}\mathbf{f}_t$ across different significance levels

| Significance | Coverage | | | |
|:---:|:---:|:---:|:---:|:---:|
| **0.1** | 0.910(0.0090) | 0.922(0.0085) | 0.943(0.0073) | 0.987(0.0036) |
| **0.05** | 0.957(0.0064) | 0.959(0.0063) | 0.974(0.0059) | 0.993(0.0026) |
| **0.01** | 0.991(0.0030) | 0.991(0.0030) | 0.994(0.0024) | 0.998(0.0014) |
| $N$ | **25** | **50** | **100** | **200** |

*Note:* This table shows the coverage rate of the first component of $\widehat{\mathbf{f}}_t - \mathbf{K}_{NT}\mathbf{f}_t$ within the confidence intervals from the theoretical distribution in Theorem 5(ii). The theoretical distribution has mean zero and variance given by the first diagonal element of $\mathbf{H}\mathbf{\Gamma}_t\mathbf{H}'$. The results are based on 1000 iterations with sample size $T = 1500$.

penalty improves the estimation accuracy, further validating the proposed method.

Finally, we compare the predictive performance of our method with QMLE by evaluating their out-of-sample forecasting accuracy under heterogeneous conditions. The forecasting error (FE) is defined as

$$\text{FE} = \left( \frac{1}{N(T - T_1)} \sum_{t=T_1+1}^{T} \|(\widehat{\mathbf{y}}_t - \mathbf{y}_t)\|_2^2 \right)^{1/2}, \tag{30}$$

where $\widehat{\mathbf{y}}_t$ denotes the predicted value using the estimated coefficients from the training sample, and $\mathbf{y}_t$ represents the actual value. Table 6 presents the forecasting error and standard deviation
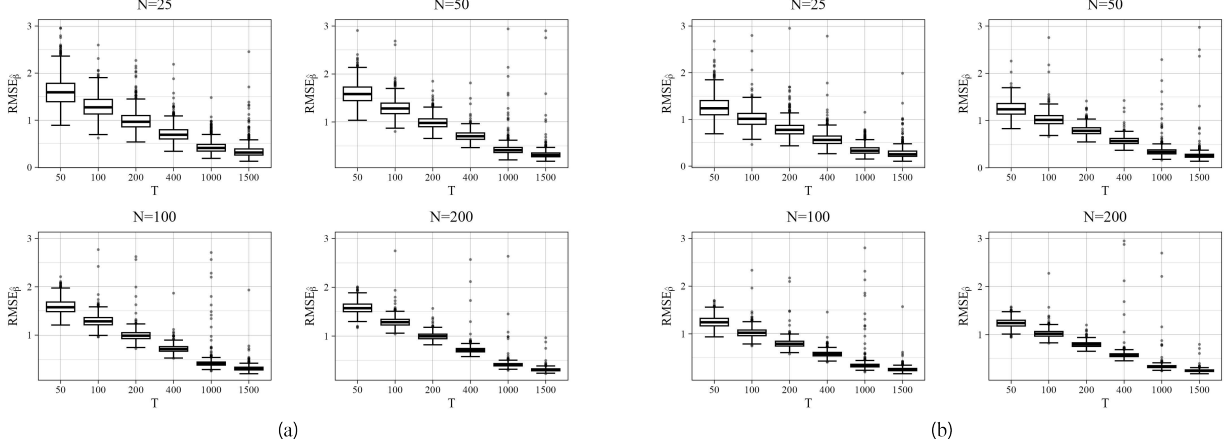
**Figure 7:** Boxplots of estimator convergence for model (18) with fixed small $\lambda_i$.
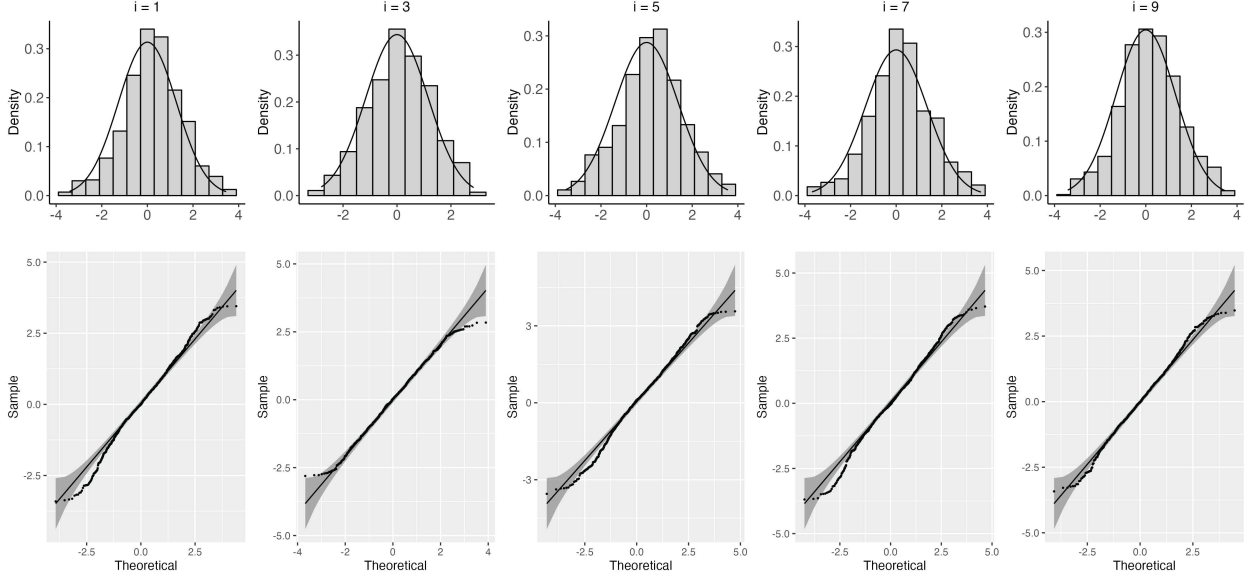


**Figure 8:** Histograms of the five spatial coefficient estimates and their corresponding empirical and theoretical distribution plots for model (18). The histograms show the distribution of the first component of $\mathbf{V}_i^H(\widetilde{\boldsymbol{\beta}}_i - \mathbf{K}_{NT}^* \boldsymbol{\beta}_i)$ as defined in Theorem 6. The superimposed normal curve represents the theoretical distribution from Theorem 6. Here, $i = 1, 3, 5, 7, 9$ correspond to the 1st, 3rd, 5th, 7th, and 9th samples in a dataset of size $N$. The results are based on 1,000 iterations with $(N, T) = (25, 3000)$.

of both method across different cross-sectional dimensions ($N$), with the out-of-sample period set from $T_1 + 1 = 321$ and $T_1 + 1 = 400$. The proposed model with lagged factor instruments ($k = 1$) achieves lower forecast error across all $N$ dimensions, outperforming QMLE. This results aligns with prior simulations that emphasize the benefits of combining shrinkage techniques with lagging factor integration to enhance accuracy.

**Table 4:** Coverage performance of $\mathbf{V}_i^H(\widetilde{\boldsymbol{\beta}}_i - \mathbf{K}_{NT}^*\boldsymbol{\beta}_i)$ under latent factor estimation across different significance levels.

| Significance | Coverage | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **0.1** | 0.343(0.016) | 0.438(0.017) | 0.531(0.017) | 0.663(0.016) | 0.819(0.013) | 0.932(0.008) | 0.979(0.005) | 1.000(0.000) |
| **0.05** | 0.399(0.016) | 0.507(0.017) | 0.595(0.017) | 0.718(0.015) | 0.866(0.011) | 0.974(0.005) | 1.000(0.000) | 1.000(0.000) |
| **0.01** | 0.477(0.017) | 0.595(0.017) | 0.692(0.016) | 0.786(0.014) | 0.946(0.008) | 1.000(0.000) | 1.000(0.000) | 1.000(0.000) |
| $T$ | **50** | **100** | **200** | **400** | **1000** | **2000** | **3000** | **5000** |

*Note:* This table shows the coverage rate of the first component of $\mathbf{V}_1^H(\widetilde{\boldsymbol{\beta}}_1 - \mathbf{K}_{NT}^*\boldsymbol{\beta}_1)$, illustrating the asymptotic performance of $\widetilde{\rho}_1$ within the confidence intervals derived from the theoretical distribution, based on 1,000 iterations with $N = 25$. The theoretical distribution has zero mean and variance corresponding to the first diagonal element of $\mathbf{X}_1^{H\prime}\mathbf{U}_1^H\mathbf{X}_1^H$, which is defined in Theorem 6.

# 6 Empirical Studies

In this section, we apply the proposed method to two arbitrage pricing case studies. The first one focuses on modeling and forecasting stock returns of the S&P 500 constituents, while the second examines quarterly changes in real housing prices across U.S. Metropolitan Statistical Areas (MSAs). For each case, we use the `R` software with a fixed random seed (`1234`) to randomly select a subset of cross-sectional units (denoted by $N$) and a subsample from the beginning of the full time span (denoted by $T$). The first 80% of each subsample is used as the training set, and the remaining 20% as the testing set for evaluating out-of-sample forecasting performance. We compare the proposed approach with the QMLE method and the classical Fama-French factor model without spatial interaction.

## 6.1 Empirical Application to Stock Returns

Companies located in close geographic proximity often share exposure to regional policies and industrial clusters, highlighting the importance of accounting for spatial dependencies. In this example, we investigate how firm locations influence stock returns. Our analysis focuses on the daily log excess returns of S&P 500 constituents from January 2004 to December 2016, comprising 3,273 time points across 205 companies. The companies are selected in the order provided by the original dataset, which is publicly available at `https://mpelger.people.stanford.edu/data-and-code`, with further details documented in Pelger (2020). As common factors, we incorporate the Fama-French variables, including the market factor (MKT), size factor (SMB), and value factor (HML). To capture spatial dependence in each selected subsample with cross-sectional dimension $N$ and time dimension $T$, we construct a spatial weight matrix following a standard geographic approach from spatial econometrics. Specifically, we define an $N \times N$ spatial weight matrix $W$, where each off-diagonal entry $w_{ij}$ is given by $w_{ij} = (s_i d_{ij})^{-1}$ for $i \neq j$, and $w_{ii} = 0$. Here, $d_{ij}$ denotes the Haversine distance between the headquarters of companies $i$ and $j$, and $s_i = \sum_{j=1}^{N} d_{ij}^{-1}$ serves as a normalization factor to ensure that each row of $W$ sums to one. This construction ensures that the

**Table 5:** Comparison of coefficient error (CE) for ridge regression estimators with latent factors across different penalized parameter $\lambda_i$ and the lagging factor impact.

| N | T | $k=1$ $(\lambda_i=10^{-9})$ | | $k=1$ $(\lambda_i=1)$ | | $k=0$ $(\lambda_i=1)$ | |
|---|---|---|---|---|---|---|---|
| | | $CE_{\widehat{\beta}}$ | $CE_{\widehat{\rho}}$ | $CE_{\widehat{\beta}}$ | $CE_{\widehat{\rho}}$ | $CE_{\widehat{\beta}}$ | $CE_{\widehat{\rho}}$ |
| 25 | 50 | 2.417(0.139) | 1.201(0.184) | 2.336(0.122) | 1.072(0.235) | 2.339(0.142) | 1.080(0.251) |
| | 100 | 2.410(0.152) | 1.194(0.202) | 2.322(0.130) | 1.068(0.247) | 2.330(0.143) | 1.066(0.263) |
| | 200 | 2.394(0.166) | 1.190(0.178) | 2.324(0.128) | 1.076(0.236) | 2.327(0.141) | 1.077(0.242) |
| | 400 | 2.396(0.178) | 1.173(0.183) | 2.314(0.128) | 1.073(0.232) | 2.336(0.131) | 1.090(0.241) |
| | 1000 | 2.419(0.144) | 1.199(0.173) | 2.310(0.143) | 1.070(0.244) | 2.334(0.139) | 1.080(0.260) |
| | 1500 | 2.435(0.153) | 1.195(0.190) | 2.323(0.126) | 1.078(0.234) | 2.323(0.139) | 1.080(0.247) |
| 50 | 50 | 2.346(0.142) | 1.001(0.230) | 2.262(0.150) | 0.838(0.306) | 2.288(0.117) | 0.973(0.277) |
| | 100 | 2.334(0.119) | 1.062(0.179) | 2.267(0.134) | 0.937(0.270) | 2.302(0.104) | 1.041(0.232) |
| | 200 | 2.325(0.100) | 1.055(0.175) | 2.259(0.134) | 0.905(0.281) | 2.282(0.105) | 0.998(0.268) |
| | 400 | 2.362(0.123) | 1.087(0.162) | 2.288(0.134) | 0.937(0.263) | 2.306(0.114) | 1.013(0.254) |
| | 1000 | 2.359(0.106) | 1.084(0.167) | 2.285(0.122) | 0.920(0.254) | 2.300(0.104) | 1.027(0.239) |
| | 1500 | 2.362(0.108) | 1.087(0.163) | 2.278(0.136) | 0.926(0.262) | 2.288(0.115) | 1.029(0.225) |
| 100 | 50 | 2.244(0.122) | 0.828(0.243) | 2.219(0.135) | 0.797(0.291) | 2.289(0.108) | 1.057(0.258) |
| | 100 | 2.261(0.102) | 0.912(0.232) | 2.227(0.119) | 0.767(0.302) | 2.299(0.097) | 1.033(0.266) |
| | 200 | 2.248(0.126) | 0.931(0.237) | 2.220(0.127) | 0.835(0.289) | 2.317(0.075) | 1.115(0.179) |
| | 400 | 2.265(0.113) | 0.947(0.223) | 2.245(0.128) | 0.880(0.304) | 2.309(0.094) | 1.081(0.228) |
| | 1000 | 2.278(0.091) | 0.965(0.194) | 2.242(0.119) | 0.929(0.266) | 2.314(0.071) | 1.139(0.185) |
| | 1500 | 2.267(0.102) | 0.939(0.214) | 2.246(0.124) | 0.870(0.291) | 2.308(0.092) | 1.089(0.208) |
| 200 | 50 | 2.230(0.116) | 0.864(0.282) | 2.230(0.086) | 0.886(0.230) | 2.352(0.055) | 1.213(0.148) |
| | 100 | 2.228(0.098) | 0.890(0.232) | 2.217(0.101) | 0.849(0.262) | 2.345(0.066) | 1.218(0.166) |
| | 200 | 2.247(0.096) | 0.945(0.248) | 2.230(0.098) | 0.921(0.231) | 2.350(0.047) | 1.230(0.104) |
| | 400 | 2.251(0.096) | 0.950(0.233) | 2.231(0.110) | 0.918(0.250) | 2.333(0.084) | 1.184(0.231) |
| | 1000 | 2.251(0.086) | 0.970(0.214) | 2.217(0.104) | 0.880(0.271) | 2.331(0.080) | 1.180(0.242) |
| | 1500 | 2.254(0.094) | 0.978(0.223) | 2.259(0.093) | 0.986(0.237) | 2.339(0.093) | 1.198(0.228) |

weights represent the relative geographic influence of company $j$ on company $i$.

Table 7 presents a detailed comparison of the forecasting performance of our proposed method against QMLE and the classical Fama-French factor model across various configurations and industry classifications. As shown in the table, our method consistently achieves lower forecasting errors than both QMLE and the factor model, highlighting its effectiveness and the value of incorporating spatial dependencies into the arbitrage pricing framework. In terms of computational efficiency, our method offers a substantial advantage over QMLE. For example, when $N = 200$ and $T = 1000$, QMLE requires over six hours on a standard CPU, while our approach produces comparable results in just a few minutes. This notable efficiency gain makes our method particularly attractive for

**Table 6:** Out-of-sample simulation evaluation of different models with best ones in boldface.

| $N$ | Proposed Model ($\lambda_i = 10^{-3}$) | | Proposed Model ($\lambda_i = 10^{-9}$) | | QML |
| | $k = 0$ | $k = 1$ | $k = 0$ | $k = 1$ | |
|---|---|---|---|---|---|
| 25 | 2.672(0.714) | **1.046(0.087)** | 2.605(0.689) | 1.047(0.088) | 1.900(0.399) |
| 50 | 2.750(1.043) | **1.047(0.090)** | 2.618(0.978) | 1.048(0.091) | 1.577(0.353) |
| 100 | 2.925(1.589) | **1.051(0.093)** | 2.634(1.386) | 1.055(0.096) | 1.357(0.224) |
| 200 | 3.423(2.753) | **1.055(0.101)** | 2.666(1.996) | 1.062(0.111) | 1.353(0.347) |

*Note:* This table compares out-of-sample forecast errors between the proposed model and QML method under varying configurations of regularization parameters ($\lambda_i = 10^{-3}$ and $\lambda_i = 10^{-9}$) and lagging factor instruments ($k = 0$ and $k = 1$). The settings for the simulation include $T = 400$ (time periods) and $K = 3$ (factor dimensions).

large-scale applications, offering a favorable trade-off between accuracy and computational cost.

**Table 7:** Forecast error comparison for stock returns with observed factors.

| Method | Proposed method | | | QMLE | | | Factor model | | |
| | $N = 100$ | $N = 150$ | $N = 200$ | $N = 100$ | $N = 150$ | $N = 200$ | $N = 100$ | $N = 150$ | $N = 200$ |
|---|---|---|---|---|---|---|---|---|---|
| $T = 500$ | **0.8601 (0.2173)** | **0.8546 (0.2396)** | **0.8462 (0.2510)** | 0.8665 (0.2177) | 0.8768 (0.2461) | 0.8772 (0.2653) | 0.8606 (0.2173) | 0.8551 (0.2396) | 0.8476 (0.2518) |
| $T = 1000$ | **0.9640 (0.1999)** | **0.9668 (0.2352)** | **0.9866 (0.2791)** | 0.9861 (0.2038) | 1.0161 (0.2505) | 1.0542 (0.3073) | 0.9646 (0.2005) | 0.9678 (0.2360) | 0.9873 (0.2801) |
| $T = 2000$ | **0.5630 (0.0641)** | **0.5704 (0.0745)** | **0.5623 (0.0818)** | 0.6042 (0.0844) | 0.6439 (0.1156) | 0.6723 (0.1510) | 0.5705 (0.0725) | 0.5821 (0.0863) | 0.5774 (0.0957) |

| GICS Class | Information Technology ($N = 36$) | | | Financials ($N = 31$) | | | Consumer Staples ($N = 19$) | | |
| | Proposed method | QMLE | Factor model | Proposed method | QMLE | Factor model | Proposed method | QMLE | Factor model |
|---|---|---|---|---|---|---|---|---|---|
| $T = 500$ | **0.7683(0.1371)** | 0.8061(0.1431) | 0.7701(0.1369) | **0.7661(0.1223)** | 0.7730(0.1256) | 0.7669(0.1223) | **0.8765(0.1396)** | 0.8766(0.1399) | 0.8767(0.1398) |
| $T = 1000$ | **0.8338(0.1252)** | 0.8697(0.1281) | 0.8352(0.1249) | **1.1530(0.2088)** | 1.1535(0.2069) | 1.1628(0.2123) | **0.9823(0.1024)** | 0.9844(0.1021) | 0.9837(0.1021) |
| $T = 2000$ | **0.6440(0.0726)** | 0.6951(0.0804) | 0.6441(0.0726) | **0.3859(0.0351)** | 0.3963(0.0366) | 0.3884(0.0351) | **0.6756(0.0556)** | 0.6821(0.0565) | 0.6793(0.0561) |

*Note:* This table compares the forecast errors of the proposed model, QMLE, and the factor model across different combinations of $N$ and $T$, as well as three industry classifications: Information Technology, Financials, and Consumer Staples. These classifications are based on the Global Industry Classification Standard (GICS).

Furthermore, we evaluate the forecasting errors for stock returns driven by unobserved factors, with the results summarized in Table 8. As shown in the table, the findings are consistent with our earlier results, further demonstrating the scalability and practicality of the proposed method for large-scale applications.

## 6.2 Empirical Application to U.S. Housing Market

Our second application examines quarterly changes in real housing prices across 377 U.S. Metropolitan Statistical Areas (MSAs) from 1975-Q1 to 2014-Q4, as studied in Aquaro et al. (2021). Due to shared supply and demand dynamics among neighboring regions, spatial models are essential for capturing such dependencies and enhancing predictive accuracy.

To account for broader economic influences—particularly the impact of stock market movements

**Table 8:** Forecast error comparison for stock returns with latent factors

| Method | Proposed method | | | QMLE | | | Factor model | | |
|---|---|---|---|---|---|---|---|---|---|
| | $N=100$ | $N=150$ | $N=200$ | $N=100$ | $N=150$ | $N=200$ | $N=100$ | $N=150$ | $N=200$ |
| $T=500$ | **0.8218(0.2138)** | **0.8251(0.2353)** | **0.8157(0.2449)** | 0.8468 (0.2132) | 0.8422 (0.2361) | 0.8353 (0.2491) | 0.8468(0.2137) | 0.8267(0.2346) | 0.8164(0.2444) |
| $T=1000$ | **0.9092(0.1919)** | **0.9315(0.2261)** | **0.9250(0.2601)** | 0.9369 (0.1968) | 0.9472 (0.2320) | 0.9605 (0.2669) | 0.9244(0.1931) | 0.9322(0.2257) | 0.9429(0.2591) |
| $T=2000$ | **0.5441(0.0690)** | **0.5501(0.0810)** | **0.5491(0.0902)** | 0.5560 (0.0704) | 0.5712 (0.0839) | 0.5698 (0.0938) | 0.5460(0.0695) | 0.5559(0.0823) | 0.5532(0.0913) |

*Note:* This table compares the forecast errors of the proposed model with QMLE and factor model under different combinations of $N$ and $T$. The number of latent factors is determined using the information criterion proposed by Bai and Ng (2002).

on real estate investment sentiment and capital allocation—we incorporate factor proxies from the previous example. For spatial dependence, we adopt the spatial weight matrix $W_{75}$ proposed in Aquaro et al. (2021), in which MSAs within a specified radius $d$ are treated as neighbors (assigned a weight of 1), while non-neighbors receive a weight of 0. The resulting matrix is row-normalized to obtain the final weight matrix $W$.

Table 9 presents a detailed comparison of forecasting errors for our proposed method, the QMLE approach from Aquaro et al. (2021), and the Fama-French factor model without spatial interactions. As shown in the table, our method consistently yields lower forecast errors in most cases, demonstrating its robustness and efficiency across different data settings and reinforcing its applicability to spatial econometric forecasting.

**Table 9:** Forecast error comparison for U.S. housing prices with observed factors.

| Method | Proposed method | | | QMLE | | | Factor model | | |
|---|---|---|---|---|---|---|---|---|---|
| | N = 20 | N=50 | N = 200 | N = 20 | N = 50 | N = 200 | N = 20 | N=50 | N = 200 |
| T = 50 | **1.2520(0.3110)** | **1.6515(0.8584)** | **1.5533(2.3174)** | 1.2680(0.2254) | 1.6605(0.7631) | 1.5863(1.5341) | 1.2881(0.2231) | 1.6931(0.7788) | 1.7533(1.5479) |
| T = 100 | **2.2001(0.6841)** | **2.1557(1.0271)** | **2.1723(2.4108)** | 2.2015(0.5007) | 2.1577(0.7837) | 2.1737(2.3700) | 2.2163(0.4752) | 2.1715(0.7257) | 2.2269(1.4651) |
| T = 150 | **2.8807(0.3411)** | **2.8876(0.4750)** | **2.7555(1.5827)** | 2.8854(0.3358) | 2.8910(0.4665) | 2.7887(1.0758) | 2.8870(0.3356) | 2.8951(0.4645) | 2.8076(1.0340) |

*Note:* This table compares the forecast errors of the proposed model, QMLE and factor model in predicting U.S. housing prices, utilizing factors from the Fama-French three-factor model and a spatial weight matrix based on geometric distances.

In summary, the comparative analyses of S&P 500 stock returns and U.S. housing prices demonstrate that our proposed method delivers superior predictive accuracy and computational efficiency, confirming its effectiveness across diverse spatial and temporal settings.

# 7 Conclusion

This paper introduced a Spatial Arbitrage Pricing Theory (SAPT) model that integrates spatial interactions with multifactor structures involving both observable and latent variables. The SAPT framework offers two key conceptual innovations for asset pricing: (1) it introduces a spatial rho parameter, serving as a counterpart to the market beta in the classical CAPM; and (2) it captures spatial correlations typically unaccounted for in traditional Arbitrage Pricing Theory (APT)

models, thereby extending the scope of standard CAPM and enhancing econometric tools for asset pricing analysis. For estimation, we proposed a generalized shrinkage Yule-Walker method that accommodates both observable and latent factors. The proposed methodology provides a flexible and computationally efficient framework for theoretical advancement and empirical research in financial and economic modeling.

# References

AHN, S. C., AND A. R. HORENSTEIN (2013): "Eigenvalue ratio test for the number of factors," *Econometrica*, 81(3), 1203–1227.

ANSELIN, L. (1988): *Spatial econometrics: methods and models*, vol. 4. Springer Science & Business Media.

AQUARO, M., N. BAILEY, AND M. H. PESARAN (2021): "Estimation and inference for spatial models with heterogeneous coefficients: an application to US house prices," *Journal of Applied Econometrics*, 36(1), 18–44.

BAI, J. (2003): "Inferential theory for factor models of large dimensions," *Econometrica*, 71(1), 135–171.

BAI, J., AND K. LI (2021): "Dynamic spatial panel data models with common shocks," *Journal of Econometrics*, 224(1), 134–160.

BAI, J., AND S. NG (2002): "Determining the number of factors in approximate factor models," *Econometrica*, 70(1), 191–221.

CARRASCO, M., AND G. TCHUENTE (2015): "Regularized LIML for many instruments," *Journal of Econometrics*, 186(2), 427–442.

CASE, A. C., H. S. ROSEN, AND J. R. HINES (1993): "Budget spillovers and fiscal policy interdependence: Evidence from the states," *Journal of Public Economics*, 52(3), 285–307.

CLIFF, A., AND J. ORD (1973): *Spatial autocorrelation.* Sage Publications Sage CA: Thousand Oaks, CA.

COCHRANE, J. (2009): *Asset pricing: Revised edition.* Princeton university press.

CRESSIE, N. (2015): *Statistics for spatial data.* John Wiley & Sons.

FAMA, E. F., AND K. R. FRENCH (1993): "Common risk factors in the returns on stocks and bonds," *Journal of Financial Economics*, 33(1), 3–56.

——— (2015): "A five-factor asset pricing model," *Journal of Financial Economics*, 116(1), 1–22.

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(4), 603–680.

FAN, J., AND Q. YAO (2003): *Nonlinear time series: nonparametric and parametric methods*, vol. 20. Springer.

FENG, G., S. GIGLIO, AND D. XIU (2020): "Taming the factor zoo: A test of new factors," *The Journal of Finance*, 75(3), 1327–1370.

FORNI, M., M. HALLIN, M. LIPPI, AND L. REICHLIN (2000): "The generalized dynamic-factor model: Identification and estimation," *Review of Economics and statistics*, 82(4), 540–554.

——— (2005): "The generalized dynamic factor model: one-sided estimation and forecasting," *Journal of the American Statistical Association*, 100(471), 830–840.

GAO, Z., Y. MA, H. WANG, AND Q. YAO (2019): "Banded spatio-temporal autoregressions," *Journal of Econometrics*, 208(1), 211–230.

GAO, Z., AND R. S. TSAY (2019): "A structural-factor approach to modeling high-dimensional time series and space-time data," *Journal of Time Series Analysis*, 40(3), 343–362.

——— (2021): "A Two-Way Transformed Factor Model for Matrix-Variate Time Series," *Econometrics and Statistics.*

——— (2022): "Modeling high-dimensional time series: A factor model with dynamically dependent factors and diverging eigenvalues," *Journal of the American Statistical Association*, 117(539), 1398–1414.

——— (2023): "Divide-and-conquer: a distributed hierarchical factor approach to modeling large-scale time series data," *Journal of the American Statistical Association*, 118(544), 2698–2711.

——— (2024): "Supervised dynamic pca: Linear dynamic forecasting with many predictors," *Journal of the American Statistical Association*, pp. 1–15.

GIGLIO, S., D. XIU, AND D. ZHANG (2025): "Test assets and weak factors," *The Journal of Finance*, 80(1), 259–319.

HANSEN, L. P. (1982): "Large sample properties of generalized method of moments estimators," *Econometrica: Journal of the econometric society*, pp. 1029–1054.

HU, J., H. DING, AND X. LIU (2023): "Arbitrage pricing with heterogeneous spatial effects and heteroscedastic disturbances," *Journal of Financial Econometrics*, 21(4), 1169–1195.

KOU, S., X. PENG, AND H. ZHONG (2018): "Asset pricing with spatial interaction," *Management Science*, 64(5), 2083–2101.

LAM, C., AND Q. YAO (2012): "Factor modeling for high-dimensional time series: inference for the number of factors," *The Annals of Statistics*, pp. 694–726.

LAM, C., Q. YAO, AND N. BATHIA (2011): "Estimation of latent factors for high-dimensional time series," *Biometrika*, 98(4), 901–918.

LEE, L.-F. (2004): "Asymptotic distributions of quasi-maximum likelihood estimators for spatial autoregressive models," *Econometrica*, 72(6), 1899–1925.

LEE, L.-F., AND J. YU (2010): "Some recent developments in spatial panel data models," *Regional Science and Urban Economics*, 40(5), 255–271.

LETTAU, M., AND M. PELGER (2020): "Estimating Latent Asset Pricing Factors," *Journal of Econometrics*, 218(1), 1–31.

LIAO, Z. (2013): "Adaptive GMM shrinkage estimation with consistent moment selection," *Econometric Theory*, 29(5), 857–904.

LIN, X., AND L.-F. LEE (2010): "GMM estimation of spatial autoregressive models with unknown heteroskedasticity," *Journal of Econometrics*, 157(1), 34–52.

LIU, X., J. GUERARD, R. CHEN, AND R. TSAY (2025): "Improving estimation of portfolio risk using new statistical factors," *Annals of Operations Research*, 346, 245–261.

MARKOWITZ, H. M. (1952): "Portfolio selection," *Journal of Finance*, 7(1), 71–91.

ONATSKI, A. (2010): "Determining the number of factors from empirical distribution of eigenvalues," *The Review of Economics and Statistics*, 92(4), 1004–1016.

PELGER, M. (2020): "Understanding systematic risk: A high-frequency approach," *The Journal of Finance*, 75(4), 2179–2220.

PESARAN, M. H., AND E. TOSETTI (2011): "Large panels with common factors and spatial correlation," *Journal of Econometrics*, 161(2), 182–202.

PIRINSKY, C., AND Q. WANG (2006): "Does corporate headquarters location matter for stock returns?," *The Journal of Finance*, 61(4), 1991–2015.

ROSS, S. A. (1976): "The arbitrage theory of capital asset pricing," *Journal of Economic Theory*, 13(3), 341–360.

SHARPE, W. F. (1964): "Capital asset prices: A theory of market equilibrium under conditions of risk," *The journal of finance*, 19(3), 425–442.

STOCK, J. H., AND M. WATSON (2002): "Macroeconomic Forecasting Using Diffusion Indexes," *Journal of Business & Economic Statistics*, 20, 147–162.

VERSHYNIN, R. (2018): *High-dimensional probability: An introduction with applications in data science*, vol. 47. Cambridge university press.

WANG, X., AND A. SHOJAIE (2021): "Joint estimation and inference for multi-experiment networks of high-dimensional point processes," *arXiv preprint arXiv:2109.11634*.

YANG, C. F. (2021): "Common factors and spatial dependence: An application to US house prices," *Econometric Reviews*, 40(1), 14–50.

YU, J., R. DE JONG, AND L.-F. LEE (2008): "Quasi-maximum likelihood estimators for spatial dynamic panel data with fixed effects when both n and T are large," *Journal of Econometrics*, 146(1), 118–134.

# Online Appendix for
# High-Dimensional Spatial Arbitrage Pricing Theory with Heterogeneous Interactions

**Abstract**

The online appendix collects the mathematical proofs that support the main text.

## IA.A    Proofs of the Theorems

We will use $C$ or $c$ to denote a generic constant the value of which may change at different places.

**Proof of Theorem 1.** We only consider a one-period economy and omit the subscript index $t$. We consider a small perturbation of the tangency portfolio $r_{j,M}$ and start with a portfolio consisting of $r_j$, $r_{j,M}$, and $r_f$ with weights $\alpha$, 1, and $-\alpha$, respective. The total wealth is $(\alpha + 1 - \alpha) = 1$. Denote the new portfolio by $r_\alpha$ and it can be written as

$$r_\alpha = r_{j,M} + \alpha r_j - \alpha r_f.$$

The variance of $r_\alpha$ is

$$\sigma_\alpha^2 = \text{Var}(r_{j,M} + \alpha r_j - \alpha r_f) = \sigma_{j,M}^2 + 2\alpha\gamma_{j,M} + \alpha^2\sigma_j^2,$$

where $\sigma_{j,M}^2 = \text{Var}(r_{j,M})$, $\sigma_j^2 = \text{Var}(r_j)$, and $\gamma_{j,M} = \text{Cov}(r_{j,M}, r_j)$. The expected return of $r_\alpha$ is

$$\mu_\alpha = \mu_{j,M} + \alpha\mu_j - \alpha r_f.$$

It follows that

$$\frac{\partial \mu_\alpha}{\partial \alpha} = \mu_j - r_f,$$

and

$$\frac{\partial \sigma_\alpha}{\partial \alpha} = \frac{1}{2}(\sigma_{j,M}^2 + 2\alpha\gamma_{j,M} + \alpha^2\sigma_j^2)^{-1/2}(2\gamma_{j,M} + 2\alpha\sigma_j^2).$$

At the tangency portfolio with $\alpha = 0$, it is known from the mean-variance theory that the slop of the capital allocation line (CAL) in Figure 1 is $\frac{\mu_{j,M} - r_f}{\sigma_{j,M}}$ as mentioned above. On the other hand

$$\frac{\partial \mu_\alpha / \partial \alpha}{\partial \sigma_\alpha / \partial \alpha}\Big|_{\alpha=0} = \frac{\mu_j - r_f}{\gamma_{j,M}/\sigma_{j,M}}.$$

From the mean-variance theory and the efficiency of the tangency portfolio $r_{j,M}$ on the frontier together with the complete market assumption as Definition 1 and Definition 2, we can conclude that the ratio between the partial derivatives above is equal to the slope of the capital allocation line:

$$\frac{\mu_{j,M} - r_f}{\sigma_{j,M}} = \frac{\mu_j - r_f}{\gamma_{j,M}/\sigma_{j,M}},$$

implying that

$$\mu_j - r_f = \frac{\gamma_{j,M}}{\sigma_{j,M}^2}(\mu_{j,M} - r_f) = \frac{\text{Cov}(r_j, \mathbf{w}_j'\mathbf{r})}{\text{Var}(\mathbf{w}_j'\mathbf{r})}(\mu_{j,M} - r_f) = \rho_j(\mu_{j,M} - r_f),$$

1

where

$$\rho_j = \frac{\mathrm{Cov}(r_j, \mathbf{w}_j'\mathbf{r})}{\mathrm{Var}(\mathbf{w}_j'\mathbf{r})},$$

is the spatial rho associated with the $j$-th asset $r_j$. This completes the proof. $\square$

**Proof of Theorem 2.** By (6), it follows that

$$\widehat{\boldsymbol{\Sigma}}_{yf} = \frac{1}{T}\sum_{t=1}^{T}\mathbf{y}_t\mathbf{f}_t' = \mathbf{D}(\boldsymbol{\rho})\mathbf{W}\frac{1}{T}\sum_{t=1}^{T}\mathbf{y}_t\mathbf{f}_t' + \mathbf{B}\frac{1}{T}\sum_{t=1}^{T}\mathbf{f}_t\mathbf{f}_t' + \frac{1}{T}\sum_{t=1}^{T}\boldsymbol{\varepsilon}_t\mathbf{f}_t'$$

and

$$\widehat{\boldsymbol{\Sigma}}_{yf}(1) = \frac{1}{T}\sum_{t=2}^{T}\mathbf{y}_t\mathbf{f}_{t-1}' = \mathbf{D}(\boldsymbol{\rho})\mathbf{W}\frac{1}{T}\sum_{t=2}^{T}\mathbf{y}_t\mathbf{f}_{t-k}' + \mathbf{B}\frac{1}{T}\sum_{t=2}^{T}\mathbf{f}_t\mathbf{f}_{t-1}' + \frac{1}{T}\sum_{t=2}^{T}\boldsymbol{\varepsilon}_t\mathbf{f}_{t-1}'.$$

Then,

$$\widehat{\boldsymbol{\Sigma}}_{yf}'\mathbf{e}_i = \widehat{\boldsymbol{\Sigma}}_{yf}'\mathbf{w}_i\rho_i + \widehat{\boldsymbol{\Sigma}}_f\mathbf{b}_i + \widehat{\boldsymbol{\Sigma}}_{\varepsilon f}'\mathbf{e}_i,$$

and

$$\widehat{\boldsymbol{\Sigma}}_{yf}(1)'\mathbf{e}_i = \widehat{\boldsymbol{\Sigma}}_{yf}(1)'\mathbf{w}_i\rho_i + \widehat{\boldsymbol{\Sigma}}_f(1)'\mathbf{b}_i + \widehat{\boldsymbol{\Sigma}}_{\varepsilon f}(1)'\mathbf{e}_i.$$

Therefore, we obtain that

$$\widehat{\mathbf{Y}}_i = \begin{pmatrix} \widehat{\boldsymbol{\Sigma}}_{yf}'\mathbf{w}_i & \widehat{\boldsymbol{\Sigma}}_f' \\ \widehat{\boldsymbol{\Sigma}}_{yf}(1)'\mathbf{w}_i & \widehat{\boldsymbol{\Sigma}}_f(1)' \end{pmatrix}\boldsymbol{\beta}_i + \begin{pmatrix} \widehat{\boldsymbol{\Sigma}}_{\varepsilon f}'\mathbf{e}_i \\ \widehat{\boldsymbol{\Sigma}}_{\varepsilon f}(1)'\mathbf{e}_i \end{pmatrix} = \widehat{\mathbf{X}}_i\boldsymbol{\beta}_i + \begin{pmatrix} \widehat{\boldsymbol{\Sigma}}_{\varepsilon f}'\mathbf{e}_i \\ \widehat{\boldsymbol{\Sigma}}_{\varepsilon f}(1)'\mathbf{e}_i \end{pmatrix}.$$

It follows that

$$\widehat{\boldsymbol{\beta}}_i(\lambda_i) = (\widehat{\mathbf{X}}_i(\lambda_i))^{-1}\widehat{\mathbf{X}}_i'\widehat{\mathbf{Y}}_i = (\widehat{\mathbf{X}}_i(\lambda_i))^{-1}\widehat{\mathbf{X}}_i'\widehat{\mathbf{X}}_i\boldsymbol{\beta}_i + (\widehat{\mathbf{X}}_i(\lambda_i))^{-1}\widehat{\mathbf{X}}_i'\begin{pmatrix} \widehat{\boldsymbol{\Sigma}}_{\varepsilon f}'\mathbf{e}_i \\ \widehat{\boldsymbol{\Sigma}}_{\varepsilon f}(1)'\mathbf{e}_i \end{pmatrix}, \qquad \text{(IA.1)}$$

where $\widehat{\mathbf{X}}_i(\lambda_i) = \widehat{\mathbf{X}}_i'\widehat{\mathbf{X}}_i + \lambda_i\mathbf{I}_{K+1}$. By Assumption 1 and a similar argument as that in (A.2) of the supplement of Gao and Tsay (2022), we can show that

$$\|\widehat{\boldsymbol{\Sigma}}_{yf}'(k)\mathbf{w}_i - \boldsymbol{\Sigma}_{yf}'(k)\mathbf{w}_i\|_F = O_p(\sqrt{\frac{N}{T}}),$$

and

$$\|\widehat{\boldsymbol{\Sigma}}_f(k) - \boldsymbol{\Sigma}_f(k)\|_F = O_p(\sqrt{\frac{1}{T}}),$$

where the first rate $\sqrt{N/T}$ can be reduced if some weak cross-sectional dependence is imposed. Furthermore, by a similar argument, we can show that

$$\widehat{\boldsymbol{\Sigma}}_{\varepsilon f}(k)'\mathbf{e}_i = \frac{1}{T}\sum_{t=k+1}^{T}\mathbf{f}_{t-k}\varepsilon_{i,t} = O_p(\sqrt{\frac{1}{T}}).$$

Therefore, if $N = o(T)$, by Assumption 6, we have

$$\|\widehat{\boldsymbol{\beta}}_i(\lambda_i) - \widehat{\mathbf{X}}_i(\lambda_i)^{-1}\widehat{\mathbf{X}}_i'\widehat{\mathbf{X}}_i\boldsymbol{\beta}_i\|_2 \leq C\|(\widehat{\mathbf{X}}_i(\lambda_i))^{-1}\widehat{\mathbf{X}}_i'\|_2 \left\| \begin{pmatrix} \widehat{\boldsymbol{\Sigma}}_{\varepsilon f}'\mathbf{e}_i \\ \widehat{\boldsymbol{\Sigma}}_{\varepsilon f}(1)'\mathbf{e}_i \end{pmatrix} \right\|_2 = O_p(T^{-1/2}),$$

and letting $\lambda_i \to 0$,

$$\|(\widehat{\mathbf{X}}_i'\widehat{\mathbf{X}}_i)(\widehat{\boldsymbol{\beta}}_i - \boldsymbol{\beta}_i)\|_2 \leq C\|\widehat{\mathbf{X}}_i'\|_2 \left\| \begin{pmatrix} \widehat{\boldsymbol{\Sigma}}_{\varepsilon f}'\mathbf{e}_i \\ \widehat{\boldsymbol{\Sigma}}_{\varepsilon f}(1)'\mathbf{e}_i \end{pmatrix} \right\|_2 = O_p(T^{-1/2}).$$

This completes the proof. $\square$

**Proof of Theorem 3.** We only prove the case when $N$ is diverging in Theorem 3(ii) as the proof for (i) is similar. By (IA.1),

$$(\widehat{\mathbf{X}}_i'\widehat{\mathbf{X}}_i)(\widehat{\boldsymbol{\beta}}_i - \boldsymbol{\beta}_i) = \widehat{\mathbf{X}}_i' \begin{pmatrix} \frac{1}{T}\sum_{t=1}^T \mathbf{f}_t\varepsilon_{i,t} \\ \frac{1}{T}\sum_{t=2}^T \mathbf{f}_{t-1}\varepsilon_{i,t} \end{pmatrix}.$$

To prove Theorem 3(ii), it is sufficient to show the following two statements,

$$\widehat{\mathbf{X}}_i'\widehat{\mathbf{X}}_i \to_p \mathbf{V}_i, \tag{IA.2}$$

and

$$\sqrt{T} \begin{pmatrix} \frac{1}{T}\sum_{t=1}^T \mathbf{f}_t\varepsilon_{i,t} \\ \frac{1}{T}\sum_{t=2}^T \mathbf{f}_{t-1}\varepsilon_{i,t} \end{pmatrix} \to_d N(\mathbf{0}, \mathbf{U}_i). \tag{IA.3}$$

By a similar argument as that in the proof of Theorem 1 above, we can show that

$$\widehat{\mathbf{X}}_i \to_p \mathbf{X}_i := \begin{pmatrix} \boldsymbol{\Sigma}_{yf}'\mathbf{w}_i & \boldsymbol{\Sigma}_f' \\ \boldsymbol{\Sigma}_{yf}(1)'\mathbf{w}_i & \boldsymbol{\Sigma}_f(1)' \end{pmatrix},$$

if $N = o(T)$. Therefore, we have

$$\widehat{\mathbf{X}}_i'\widehat{\mathbf{X}}_i \to_p \begin{pmatrix} \mathbf{w}_i'\boldsymbol{\Sigma}_{yf}\boldsymbol{\Sigma}_{yf}'\mathbf{w}_i + \mathbf{w}_i'\boldsymbol{\Sigma}_{yf}(1)\boldsymbol{\Sigma}_{yf}'(1)\mathbf{w}_i & \mathbf{w}_i'\boldsymbol{\Sigma}_{yf}\boldsymbol{\Sigma}_f\mathbf{w}_i + \mathbf{w}_i'\boldsymbol{\Sigma}_{yf}(1)\boldsymbol{\Sigma}_f'(1) \\ \boldsymbol{\Sigma}_f\boldsymbol{\Sigma}_{yf}'\mathbf{w}_i + \boldsymbol{\Sigma}_f(1)\boldsymbol{\Sigma}_{yf}'(1)\mathbf{w}_i & \boldsymbol{\Sigma}_f^2 + \boldsymbol{\Sigma}_f(1)\boldsymbol{\Sigma}_f'(1) \end{pmatrix} = \mathbf{V}_i.$$

By a similar argument as that in the proof of Theorem 2, we can show that

$$\widehat{\mathbf{X}}_i'\widehat{\mathbf{X}}_i - \mathbf{X}_i'\mathbf{X}_i = O_p(N^{-1/2}T^{1/2}),$$

3

which implies (IA.2) if $N = o(T)$. To show (IA.3), it is sufficient to prove that, for any vector $\mathbf{a} = (\mathbf{a}_1', \mathbf{a}_2')'$ with $\mathbf{a}_1 \in R^K$ and $\mathbf{a}_2 \in R^K$,

$$\sqrt{T}\mathbf{a}' \left( \begin{array}{c} \frac{1}{T}\sum_{t=1}^{T}\mathbf{f}_t\varepsilon_{i,t} \\ \frac{1}{T}\sum_{t=2}^{T}\mathbf{f}_{t-1}\varepsilon_{i,t} \end{array} \right) = \mathbf{a}_1'\frac{1}{\sqrt{T}}\sum_{t=1}^{T}\mathbf{f}_t\varepsilon_{i,t} + \mathbf{a}_2'\frac{1}{\sqrt{T}}\sum_{t=1}^{T}\mathbf{f}_{t-1}\varepsilon_{i,t}$$

(IA.4)

is asymptotically normal. Define

$$\mathbf{S}_{N,T} = a_1\frac{1}{\sqrt{T}}\sum_{t=1}^{T}\mathbf{f}_t\varepsilon_{i,t} + \mathbf{a}_2'\frac{1}{\sqrt{T}}\sum_{t=2}^{T}\mathbf{f}_{t-1}\varepsilon_{i,t},$$

(IA.5)

we only need to show the asymptotic normality of $\mathbf{S}_{N,T}$. By Schwarz's Inequality and Assumptions 2 and 5, we can derive that

$$E|\mathbf{f}_t\varepsilon_{i,t}|^\gamma \leq (E|\mathbf{f}_t|^{2\gamma})^{1/2}(E|\varepsilon_{it}|^{2\gamma})^{1/2} < \infty.$$

We now calculate the variance of $\mathbf{S}_{N,T}$. Since it involves 16 terms in total and we start with the first term in (IA.5). By definition and an elementary argument, we have

$$\text{Var}(\frac{1}{\sqrt{T}}\sum_{t=1}^{T}\mathbf{f}_t\varepsilon_{i,t}) = \mathbf{\Sigma}_{f\varepsilon_i}(0) + \sum_{j=1}^{T-1}(1 - \frac{j}{T})\mathbf{\Sigma}_{f\varepsilon_i}(0,j).$$

Note that $\sum_{j=1}^{\infty}\alpha_N(j)^{1-2/\gamma} < \infty$ from Assumption 1, by Proposition 2.5 of Fan and Yao (2003), we have

$$\sup_i\sum_{j=1}^{\infty}|\mathbf{\Sigma}_{f\varepsilon_i}(0,j)| \leq C\sup_i\sum_{j=1}^{\infty}\alpha(j)^{1-2/\gamma}(E|\mathbf{f}_t|^{2\gamma})^{1/\gamma}(E|\varepsilon_{i,t}|^{2\gamma})^{1/\gamma} < \infty.$$

We can calculate all the terms of $\mathbf{S}_{N,T}$ and sum them up, by the Dominated Convergence theorem, we have

$$\text{Var}\,(\mathbf{S}_{N,T}) \to \mathbf{a}'\mathbf{U}_i\mathbf{a}.$$

To show the asymptotic normality of $\mathbf{S}_{N,T}$, we employ the small-block and large-block techniques commonly used for weakly dependent data. Specifically, we partition the set $\{1, ..., T\}$ into $2k_T + 1$ subsets with large blocks of size $l_T$, small blocks of size $s_T$, and the last remaining set of size $T - k_T(l_T + s_T)$. Let

$$l_T = [\sqrt{T}/\log(T)], \ s_T = [\sqrt{T}\log(T)]^\delta, \ k_T = [T/(l_T + s_T)],$$

where $[x]$ is the greatest integer less than or equal to $x$, and $1 - 2/\gamma \leq \delta < 1$. It is not hard to see that

$$l_T/\sqrt{T} \to 0, \ s_T/l_T \to 0, \ \text{and } k_T = O(\sqrt{T}\log(T)).$$

4

By Assumption 1 that $\sum_{j=1}^{\infty} \alpha_N(j)^{1-2/\gamma}$, we have $\alpha_N(s_T) = o(s_T^{-\gamma/(\gamma-2)})$. It follows that

$$k_T \alpha_N(s_T) = o(k_T/s_T^{\gamma/(\gamma-2)}) = o(1).$$

Then, we rewrite $\mathbf{S}_{N,T}$ as

$$\mathbf{S}_{N,T} = \mathbf{a}_1' \frac{1}{\sqrt{T}} \sum_{j=1}^{k_T} \boldsymbol{\xi}_j^{(1)} + \mathbf{a}_2' \frac{1}{\sqrt{T}} \sum_{j=1}^{k_T} \boldsymbol{\xi}_j^{(2)} + \mathbf{a}_1' \frac{1}{\sqrt{T}} \sum_{j=1}^{k_T} \boldsymbol{\eta}_j^{(1)} + \mathbf{a}_2' \frac{1}{\sqrt{T}} \sum_{j=1}^{k_T} \boldsymbol{\eta}_j^{(2)}$$
$$+ \mathbf{a}_1' \frac{1}{\sqrt{T}} \boldsymbol{\zeta}_j^{(1)} + \mathbf{a}_2' \frac{1}{\sqrt{T}} \boldsymbol{\zeta}_j^{(2)}, \tag{IA.6}$$

where

$$\boldsymbol{\xi}_j^{(1)} = \sum_{t=(j-1)(l_T+s_T)+1}^{jl_T+(j-1)s_T} \mathbf{f}_t \varepsilon_{i,t}, \quad \boldsymbol{\eta}_j^{(1)} = \sum_{t=jl_T+(j-1)s_T+1}^{j(l_T+s_T)} \mathbf{f}_t \varepsilon_{i,t},$$

$$\boldsymbol{\zeta}_j^{(1)} = \sum_{t=k_T(l_T+s_T)+1}^{T} \mathbf{f}_t \varepsilon_{i,t}, \quad \boldsymbol{\xi}_j^{(2)} = \sum_{t=(j-1)(l_T+s_T)+1}^{jl_T+(j-1)s_T} \mathbf{f}_{t-1} \varepsilon_{i,t},$$

$$\boldsymbol{\eta}_j^{(2)} = \sum_{t=jl_T+(j-1)s_T+1}^{j(l_T+s_T)} \mathbf{f}_{t-1} \varepsilon_{i,t}, \quad \boldsymbol{\zeta}_j^{(2)} = \sum_{t=k_T(l_T+s_T)+1}^{T} \mathbf{f}_{t-1} \varepsilon_{i,t},$$

Note that $\alpha_N(T) = o(T^{2/\gamma-1})$, $k_T s_T/T \to 0$, and $(l_T+s-T)/T \to 0$, it follows from Proposition 2.7 of Fan and Yao (2003) that

$$\frac{1}{\sqrt{T}} \sum_{j=1}^{k_T} \boldsymbol{\eta}_j^{(l)} = o_p(1), \text{ and } \frac{1}{\sqrt{T}} \boldsymbol{\zeta}_j^{(l)} = o_p(1), \ l = 1, 2, 3, 4.$$

Then,

$$\mathbf{S}_{N,T} = \mathbf{a}_1' \frac{1}{\sqrt{T}} \sum_{j=1}^{k_T} \boldsymbol{\xi}_j^{(1)} + \mathbf{a}_2' \frac{1}{\sqrt{T}} \sum_{j=1}^{k_T} \boldsymbol{\xi}_j^{(2)} + o_p(1).$$

By a similar argument as Theorem 2.21 of Fan and Yao (2003), we can show that

$$\sqrt{T} \mathbf{a}' \begin{pmatrix} \frac{1}{T} \sum_{t=1}^{T} \mathbf{f}_t \varepsilon_{i,t} \\ \frac{1}{T} \sum_{t=2}^{T} \mathbf{f}_{t-1} \varepsilon_{i,t} \end{pmatrix} \longrightarrow_d N(0, \mathbf{a}' \mathbf{U}_i \mathbf{a}).$$

We replace $\mathbf{a}$ by $(\mathbf{U}_i^{-1/2})' \mathbf{a}$ and obtain

$$\sqrt{T} \mathbf{a}' \mathbf{U}_i^{-1/2} \begin{pmatrix} \frac{1}{T} \sum_{t=1}^{T} \mathbf{f}_t \varepsilon_{i,t} \\ \frac{1}{T} \sum_{t=2}^{T} \mathbf{f}_{t-1} \varepsilon_{i,t} \end{pmatrix} \longrightarrow_d N(0, 1),$$

which implies that

$$\sqrt{T} \mathbf{U}_i^{-1/2} \begin{pmatrix} \frac{1}{T} \sum_{t=1}^{T} \mathbf{f}_t \varepsilon_{i,t} \\ \frac{1}{T} \sum_{t=2}^{T} \mathbf{f}_{t-1} \varepsilon_{i,t} \end{pmatrix} \longrightarrow_d N(\mathbf{0}, \mathbf{I}_{2K}). \tag{IA.7}$$

Therefore,

$$\sqrt{T}\mathbf{V}_i(\widehat{\boldsymbol{\beta}}_i - \boldsymbol{\beta}_i) \longrightarrow_d N(\mathbf{0}, \mathbf{X}_i'\mathbf{U}_i\mathbf{X}_i). \tag{IA.8}$$

Theorem 3 follows from (IA.2) and (IA.8). This completes the proof. $\square$

**Proof of Theorem 4.** Note that

$$\begin{aligned}
\widehat{\boldsymbol{\Sigma}}_y(k) &= \frac{1}{T}\sum_{t=k+1}^{T} \mathbf{y}_t\mathbf{y}_{t-k}' \\
&= \frac{1}{T}\sum_{t=k+1}^{T} \{\boldsymbol{\Lambda}\mathbf{f}_t\mathbf{f}_{t-k}'\boldsymbol{\Lambda}' + \boldsymbol{\Lambda}\mathbf{f}_t\boldsymbol{\xi}_{t-k}' + \boldsymbol{\xi}_t\mathbf{f}_{tk}'\boldsymbol{\Lambda}' + \boldsymbol{\xi}_t\boldsymbol{\xi}_{t-k}'\} \\
&= \boldsymbol{\Lambda}\frac{1}{T}\sum_{t=k+1}^{T} (\mathbf{f}_t\mathbf{f}_{t-k}'\boldsymbol{\Lambda}' + \mathbf{f}_t\boldsymbol{\xi}_{t-k}') + \frac{1}{T}\sum_{t=k+1}^{T} (\boldsymbol{\xi}_t\mathbf{f}_{t-k}'\boldsymbol{\Lambda}' + \boldsymbol{\xi}_t\boldsymbol{\xi}_{t-k}') \\
&= \boldsymbol{\Lambda}\mathbf{G}_{1,k} + \mathbf{G}_{2,k}, \tag{IA.9}
\end{aligned}$$

where

$$\mathbf{G}_{1,k} = \frac{1}{T}\sum_{t=k+1}^{T} (\mathbf{f}_t\mathbf{f}_{t-k}'\boldsymbol{\Lambda}' + \mathbf{f}_t\boldsymbol{\xi}_{t-k}'), \ \ \mathbf{G}_{2,k} = \frac{1}{T}\sum_{t=k+1}^{T} (\boldsymbol{\xi}_t\mathbf{f}_{t-k}'\boldsymbol{\Lambda}' + \boldsymbol{\xi}_t\boldsymbol{\xi}_{t-k}').$$

It follows from the definition of $\widehat{\mathbf{M}}$ in (17) that

$$\begin{aligned}
\widehat{\mathbf{M}} &= \sum_{k=1}^{k_0} \widehat{\boldsymbol{\Sigma}}_y(k)\widehat{\boldsymbol{\Sigma}}_y'(k) \\
&= \sum_{k=1}^{k_0} (\boldsymbol{\Lambda}\mathbf{G}_{1,k} + \mathbf{G}_{2,k})(\boldsymbol{\Lambda}\mathbf{G}_{1,k} + \mathbf{G}_{2,k})' \\
&= \boldsymbol{\Lambda}\sum_{k=1}^{k_0} \mathbf{G}_{1,k}\mathbf{G}_{1,k}'\boldsymbol{\Lambda}' + \sum_{k=1}^{k_0} (\boldsymbol{\Lambda}\mathbf{G}_{1,k}\mathbf{G}_{2,k}' + \mathbf{G}_{2,k}\mathbf{G}_{1,k}'\boldsymbol{\Lambda}' + \mathbf{G}_{2,k}\mathbf{G}_{2,k}'). \tag{IA.10}
\end{aligned}$$

Let $\widehat{\mathbf{V}}_{NT} \in R^r$ be a diagonal matrix with diagonal elements being the top $K$ eigenvalues of $\widehat{\mathbf{M}}$, it follows from Assumptions 3 and 4 that $\widehat{\mathbf{V}}_{NT} \asymp O(N^2)$. Since the columns of $\widehat{\boldsymbol{\Lambda}}$ are the eigenvectors of $\widehat{\mathbf{M}}$, it follows that

$$\widehat{\mathbf{M}}\widehat{\boldsymbol{\Lambda}} = \widehat{\boldsymbol{\Lambda}}\widehat{\mathbf{V}}_{NT},$$

implying that

$$\begin{aligned}
\widehat{\boldsymbol{\Lambda}} &= \widehat{\mathbf{M}}\widehat{\boldsymbol{\Lambda}}\widehat{\mathbf{V}}_{NT}^{-1} \\
&= \boldsymbol{\Lambda}\sum_{k=1}^{k_0} \mathbf{G}_{1,k}\mathbf{G}_{1,k}'\boldsymbol{\Lambda}'\widehat{\boldsymbol{\Lambda}}\widehat{\mathbf{V}}_{NT}^{-1} + \sum_{k=1}^{k_0} \left[\boldsymbol{\Lambda}\mathbf{G}_{1,k}\mathbf{G}_{2,k}' + \mathbf{G}_{2,k}\mathbf{G}_{1,k}'\boldsymbol{\Lambda}' + \mathbf{G}_{2,k}\mathbf{G}_{2,k}'\right]\widehat{\boldsymbol{\Lambda}}\widehat{\mathbf{V}}_{NT}^{-1}. \tag{IA.11}
\end{aligned}$$

Let $\mathbf{H}_{NT}' = \sum_{k=1}^{k_0} \mathbf{G}_{1,k}\mathbf{G}_{1,k}'\boldsymbol{\Lambda}'\widehat{\boldsymbol{\Lambda}}\widehat{\mathbf{V}}_{NT}^{-1}$, it follows that $\mathbf{H} = O_p(1)$ and $\mathbf{H}^{-1} = O_p(1)$. Then (IA.11)

implies that

$$\widehat{\mathbf{\Lambda}} - \mathbf{\Lambda}\mathbf{H}'_{NT} = \sum_{k=1}^{k_0} \left[ \mathbf{\Lambda}\mathbf{G}_{1,k}\mathbf{G}'_{2,k} + \mathbf{G}_{2,k}\mathbf{G}'_{1,k}\mathbf{\Lambda}' + \mathbf{G}_{2,k}\mathbf{G}'_{2,k} \right] \widehat{\mathbf{\Lambda}}\widehat{\mathbf{V}}_{NT}^{-1}. \tag{IA.12}$$

First, by Assumption 1 and a similar argument as that in (A.2) of the supplement of Gao and Tsay (2022), we can show that

$$\|\mathbf{G}_{1,k}\|_F = \|\frac{1}{T}\sum_{t=k+1}^{T}(\mathbf{f}_t\mathbf{f}'_{t-k}\mathbf{\Lambda}' + \mathbf{f}_t\boldsymbol{\xi}'_{t-k})\|_F = O_p(\sqrt{N}) + O_p(1 + \sqrt{\frac{N}{T}}) = O_p(\sqrt{N}),$$

and

$$\|\mathbf{G}_{2,k}\|_F = \|\frac{1}{T}\sum_{t=k+1}^{T}(\boldsymbol{\xi}_t\mathbf{f}'_{t-k}\mathbf{\Lambda}' + \boldsymbol{\xi}_t\boldsymbol{\xi}'_{t-k})\|_F = O_p(\sqrt{\frac{N}{T}}\sqrt{N}) + O_p(\sqrt{\frac{N^2}{T}}) = O_p(\sqrt{\frac{N^2}{T}}).$$

Then, it follows from (IA.12) and the above rates that

$$\|\widehat{\mathbf{\Lambda}} - \mathbf{\Lambda}\mathbf{H}'\|_F = O_p(\sqrt{N}\sqrt{N}\sqrt{\frac{N^2}{T}} + \sqrt{N}\sqrt{N}\sqrt{\frac{N^2}{T}} + \frac{N^2}{T})O_p(\sqrt{N}/N^2) = O_p(\sqrt{\frac{N}{T}}),$$

implying that

$$\frac{1}{\sqrt{N}}\|\widehat{\mathbf{\Lambda}} - \mathbf{\Lambda}\mathbf{H}'\|_F = O_p(T^{-1/2}).$$

This completes the proof of Theorem 4. $\square$

**Lemma 1.** *Let Assumptions 1–8 hold. Then, as $N, T \to \infty$,*

$$\mathbf{H}_{NT}\mathbf{H}'_{NT} = \mathbf{I}_K + O_p(T^{-1/2}), \text{ and } \mathbf{H}'_{NT}\mathbf{H}_{NT} = \mathbf{I}_K + O_p(T^{-1/2}).$$

**Proof.** First, note that

$$\begin{aligned}
\mathbf{H}_{NT}\mathbf{H}'_{NT} - \mathbf{I}_K &= \mathbf{H}_{NT}\frac{\mathbf{\Lambda}'\mathbf{\Lambda}}{N}\mathbf{H}'_{NT} - \frac{\widehat{\mathbf{\Lambda}}'\widehat{\mathbf{\Lambda}}}{N} \\
&= \frac{1}{\sqrt{N}}(\mathbf{H}_{NT}\mathbf{\Lambda}' - \widehat{\mathbf{\Lambda}}')\frac{1}{\sqrt{N}}\mathbf{\Lambda}\mathbf{H}'_{NT} + \frac{1}{\sqrt{N}}\widehat{\mathbf{\Lambda}}'(\mathbf{\Lambda}\mathbf{H}'_{NT} - \widehat{\mathbf{\Lambda}})/\sqrt{N}. \tag{IA.13}
\end{aligned}$$

Then, it follows from Theorem 3 that

$$\|\mathbf{H}_{NT}\mathbf{H}'_{NT} - \mathbf{I}_K\|_F = O_p(T^{-1/2}).$$

Furthermore, since $\mathbf{H}_{NT} = O_p(1)$ and $\mathbf{H}_{NT}^{-1} = O_p(1)$, then

$$\mathbf{H}'_{NT}\mathbf{H}_{NT}\mathbf{H}'_{NT} = \mathbf{H}'_{NT} + O_p(T^{-1/2}),$$

we multiply $\mathbf{H}_{NT}^{-1}$ on the right of both sides and obtain

$$\mathbf{H}_{NT}'\mathbf{H}_{NT} = \mathbf{I}_K + O_p(T^{-1/2}).$$

This completes the proof. $\square$

**Lemma 2.** *Let Assumptions 1–8 hold. Then, as $N, T \to \infty$,*

$$\widehat{\mathbf{\Lambda}}'\widehat{\mathbf{M}}\widehat{\mathbf{\Lambda}} = \widehat{\mathbf{V}}_{NT} \to_p \mathbf{V},$$

*where $\mathbf{V}$ is a diagonal matrix consisting of the top $K$ eigenvalues of $\mathbf{M}$ defined in (16).*

**Proof.** The proof is similar to Theorem 1 of Lam and Yao (2012). We omit the details to save space. $\square$

**Lemma 3.** *Let Assumptions 1–8 hold. Then there exists an orthogonal matrix $\mathbf{H} \in R^K$ such that $\mathbf{H}_{NT} \to_p \mathbf{H}$ with probability tending to one as $N, T \to \infty$.*

**Proof.** Note that

$$\mathbf{G}_{1,k} = \widehat{\mathbf{\Sigma}}_f(k)\mathbf{\Lambda}' + \widehat{\mathbf{\Sigma}}_{f\xi}(k).$$

If $N = o(T)$, we have that $\widehat{\mathbf{\Sigma}}_f(k) \to_p \mathbf{\Sigma}_f(k)$ and $\widehat{\mathbf{\Sigma}}_{f\xi}(k) \to_p \mathbf{\Sigma}_{f\xi}(k)$. By definition,

$$\mathbf{H}_{NT}' = \sum_{k=1}^{k_0} \mathbf{G}_{1,k}\mathbf{G}_{1,k}'\mathbf{\Lambda}'\widehat{\mathbf{\Lambda}}\widehat{\mathbf{V}}_{NT}^{-1} = \sum_{k=1}^{k_0} \mathbf{G}_{1,k}\mathbf{G}_{1,k}'\mathbf{\Lambda}'\mathbf{\Lambda}\mathbf{H}_{NT}'\mathbf{V}^{-1} + o_p(1).$$

Then, by Lemma 1,

$$\mathbf{H}_{NT}'\frac{\mathbf{V}}{N^2}\mathbf{H}_{NT} = \frac{1}{N}\sum_{k=1}^{k_0}(\mathbf{\Sigma}_f(k)\mathbf{\Lambda}' + \mathbf{\Sigma}_{f\xi}(k)))(\mathbf{\Sigma}_f(k)\mathbf{\Lambda}' + \mathbf{\Sigma}_{f\xi}(k)))' + o_p(1)$$

$$= \sum_{k=1}^{k_0}\mathbf{\Sigma}_f(k)\mathbf{\Sigma}_f'(k) + o_p(1). \tag{IA.14}$$

Therefore, $\mathbf{H}_{NT}$ will converge to the matrix consisting of the eigenvectors of $\sum_{k=1}^{k_0}\mathbf{\Sigma}_f(k)\mathbf{\Sigma}_f'(k)$, denoted by $\mathbf{H}$. This completes the proof. $\square$

**Proof of Theorem 5.** We first show Theorem 5(i). Note that

$$\begin{aligned}
\widehat{\mathbf{f}}_t &= \frac{1}{N}\widehat{\mathbf{\Lambda}}'\mathbf{y}_t \\
&= \frac{1}{N}\widehat{\mathbf{\Lambda}}'(\mathbf{\Lambda}\mathbf{f}_t + \boldsymbol{\xi}_t) \\
&= \frac{1}{N}\widehat{\mathbf{\Lambda}}'\mathbf{\Lambda}\mathbf{f}_t + \frac{1}{N}\widehat{\mathbf{\Lambda}}'\boldsymbol{\xi}_t \\
&= \mathbf{K}_{NT}\mathbf{f}_t + \frac{1}{N}(\widehat{\mathbf{\Lambda}} - \mathbf{\Lambda}\mathbf{H}_{NT}')'\boldsymbol{\xi}_t + \frac{1}{N}\mathbf{H}_{NT}\mathbf{\Lambda}'\boldsymbol{\xi}_t,
\end{aligned} \tag{IA.15}$$

8

where $\mathbf{K}_{NT} = \frac{1}{N}\widehat{\boldsymbol{\Lambda}}'\boldsymbol{\Lambda}$ which has the same limit $\mathbf{H}$ as that of $\mathbf{H}_{NT}$, but they are not identical in finite samples. By Assumption 7, it is not hard to show that

$$\max_{1\leq t\leq t}|f_{i,t}| = O_p(\log(T)), \text{ and } \max_{1\leq t\leq T}|\varepsilon_{i,t}| = O_p(\log(T)).$$

Therefore,

$$\max_{1\leq t\leq T}\|\frac{1}{N}(\widehat{\boldsymbol{\Lambda}} - \boldsymbol{\Lambda}\mathbf{H}'_{NT})'\boldsymbol{\xi}_t\|_F \leq \frac{1}{\sqrt{N}}\|\widehat{\boldsymbol{\Lambda}} - \boldsymbol{\Lambda}\mathbf{H}'_{NT}\|_F \max_{1\leq t\leq T}\|\frac{\boldsymbol{\xi}_t}{\sqrt{N}}\|_F = O_p(T^{-1/2}\log(T)),$$

and

$$\max_{1\leq t\leq T}\|\frac{1}{N}\mathbf{H}_{NT}\boldsymbol{\Lambda}'\boldsymbol{\xi}_t\|_F \leq C\frac{1}{\sqrt{N}}\|\mathbf{H}_{NT}\|_F \max_{1\leq t\leq T}\|\frac{1}{\sqrt{N}}\boldsymbol{\Lambda}'\boldsymbol{\xi}_t\|_F = O_p(N^{-1/2}\log(T)).$$

Then, it follows from the above rates that

$$\max_{1\leq t\leq t}\|\widehat{\mathbf{f}}_t - \mathbf{K}_{NT}\mathbf{f}_t\|_F = O_p((\frac{1}{\sqrt{T}} + \frac{1}{\sqrt{N}})\log(T)).$$

This proves Theorem 5(i).

For Theorem 5(ii), if $N = o(T)$,

$$\sqrt{N}(\widehat{\mathbf{f}}_t - \mathbf{K}_{NT}\mathbf{f}_t) = \mathbf{H}\frac{1}{\sqrt{N}}\boldsymbol{\Lambda}'\boldsymbol{\xi}_t + o_p(1). \tag{IA.16}$$

By Assumption 8, we have

$$\sqrt{N}(\widehat{\mathbf{f}}_t - \mathbf{K}_{NT}\mathbf{f}_t) \longrightarrow_d N(0, \mathbf{H}\boldsymbol{\Gamma}_t\mathbf{H}').$$

This completes the proof. $\square$

**Proof of Theorem 6.** We use $\mathbf{H}$ instead of $\mathbf{H}_{NT}$ for simplicity in this proof. By definition,

$$\begin{aligned}
\widetilde{\boldsymbol{\Sigma}}_{yf}(k) =& \frac{1}{T}\sum_{t=k+1}^{T}\mathbf{y}_t\widehat{\mathbf{f}}'_{t-k} \\
=& \frac{1}{T}\sum_{t=k+1}^{T}\mathbf{D}(\boldsymbol{\rho})\mathbf{W}\mathbf{y}_t\widehat{\mathbf{f}}'_{t-k} + \mathbf{B}\frac{1}{T}\sum_{t=k+1}^{T}\mathbf{f}_t\widehat{\mathbf{f}}'_{t-k} + \frac{1}{T}\sum_{t=k+1}^{T}\boldsymbol{\varepsilon}_t\widehat{\mathbf{f}}'_{t-k} \\
=& \mathbf{D}(\boldsymbol{\rho})\mathbf{W}\widetilde{\boldsymbol{\Sigma}}_{yf}(k) + \mathbf{B}\mathbf{K}_{NT}^{-1}\widetilde{\boldsymbol{\Sigma}}_f(k) + \mathbf{B}\frac{1}{T}\sum_{t=k+1}^{T}(\mathbf{f}_t - \mathbf{K}_{NT}^{-1}\widehat{\mathbf{f}}_t)\widehat{\mathbf{f}}'_{t-k} + \frac{1}{T}\sum_{t=k+1}^{T}\boldsymbol{\varepsilon}_t\widehat{\mathbf{f}}'_{t-k}. \tag{IA.17}
\end{aligned}$$

Then, it follows that

$$\widetilde{\boldsymbol{\Sigma}}_{yf}(k)'\mathbf{e}_i = \widetilde{\boldsymbol{\Sigma}}_{yf}(k)'\mathbf{w}_i\rho_i + \widetilde{\boldsymbol{\Sigma}}_f(k)'(\mathbf{K}'_{NT})^{-1}\mathbf{b}_i + \frac{1}{T}\sum_{t=k+1}^{T}\widehat{\mathbf{f}}_{t-k}(\mathbf{f}_t - \mathbf{K}_{NT}^{-1}\widehat{\mathbf{f}}_t)'\mathbf{b}_i + \frac{1}{T}\sum_{t=k+1}^{T}\widehat{\mathbf{f}}_{t-k}\varepsilon_{i,t}. \tag{IA.18}$$

9

We now analyze the last two terms. First,

$$\frac{1}{T}\sum_{t=k+1}^{T}\widehat{\mathbf{f}}_{t-k}(\mathbf{f}_t - \mathbf{K}_{NT}^{-1}\widehat{\mathbf{f}}_t)'\mathbf{b}_i = \frac{1}{T}\sum_{t=k+1}^{T}(\widehat{\mathbf{f}}_{t-k} - \mathbf{K}_{NT}\mathbf{f}_{t-k})(\mathbf{K}_{NT}\mathbf{f}_t - \widehat{\mathbf{f}}_t)'(\mathbf{K}'_{NT})^{-1}\mathbf{b}_i$$

$$+ \frac{1}{T}\sum_{t=k+1}^{T}\mathbf{K}_{NT}\mathbf{f}_{t-k}(\mathbf{K}_{NT}\mathbf{f}_t - \widehat{\mathbf{f}}_t)'(\mathbf{K}'_{NT})^{-1}\mathbf{b}_i, \qquad \text{(IA.19)}$$

where

$$\|\frac{1}{T}\sum_{t=k+1}^{T}(\widehat{\mathbf{f}}_{t-k} - \mathbf{K}_{NT}\mathbf{f}_{t-k})(\mathbf{K}_{NT}\mathbf{f}_t - \widehat{\mathbf{f}}_t)'(\mathbf{K}'_{NT})^{-1}\mathbf{b}_i\|_2 = O_p((\frac{1}{T} + \frac{1}{N})\log(T)^2).$$

By (IA.16), and $\mathbf{f}_{t-k}$ and $\boldsymbol{\xi}_t$ are uncorrelated, we have

$$\frac{1}{T}\sum_{t=k+1}^{T}\mathbf{K}_{NT}\mathbf{f}_{t-k}(\mathbf{K}_{NT}\mathbf{f}_t - \widehat{\mathbf{f}}_t)'(\mathbf{K}'_{NT})^{-1}\mathbf{b}_i = \frac{1}{T}\sum_{t=k+1}^{T}\mathbf{K}_{NT}\mathbf{f}_{t-k}\boldsymbol{\xi}'_t(\boldsymbol{\Lambda}\mathbf{H}'_{NT} - \widehat{\boldsymbol{\Lambda}})(\mathbf{K}_{NT})^{-1}\mathbf{b}_i/N$$

$$- \frac{1}{T}\sum_{t=k+1}^{T}\mathbf{K}_{NT}\mathbf{f}_{t-k}\boldsymbol{\xi}'_t\boldsymbol{\Lambda}\mathbf{H}'_{NT}(\mathbf{K}_{NT})^{-1}\mathbf{b}_i/N$$

$$= O_p(\frac{1}{T} + \frac{1}{\sqrt{NT}}). \qquad \text{(IA.20)}$$

Next, we consider the last term of (IA.18).

$$\frac{1}{T}\sum_{t=k+1}^{T}\widehat{\mathbf{f}}_{t-k}\varepsilon_{i,t} = \frac{1}{T}\sum_{t=k+1}^{T}\mathbf{K}_{NT}\mathbf{f}_{t-k}\varepsilon_{i,t} + \frac{1}{T}\sum_{t=k+1}^{T}(\widehat{\mathbf{f}}_{t-k} - \mathbf{K}_{NT}\mathbf{f}_{t-k})\varepsilon_{i,t}.$$

By a similar argument as (IA.20), we can show that

$$\|\frac{1}{T}\sum_{t=k+1}^{T}(\widehat{\mathbf{f}}_{t-k} - \mathbf{K}_{NT}\mathbf{f}_{t-k})\varepsilon_{i,t}\|_F = O_p(\frac{1}{T} + \frac{1}{\sqrt{NT}}).$$

Then it follows from (IA.18) that

$$\widetilde{\mathbf{Y}}_i = \widetilde{\mathbf{X}}_i\mathbf{K}^*_{NT}\boldsymbol{\beta}_i + \begin{pmatrix} \mathbf{K}_{NT}\frac{1}{T}\sum_{t=1}^{T}\mathbf{f}_t\varepsilon_{i,t} \\ \mathbf{K}_{NT}\frac{1}{T}\sum_{t=2}^{T}\mathbf{f}_{t-1}\varepsilon_{i,t} \end{pmatrix} + O_p(\frac{1}{\sqrt{NT}} + (\frac{1}{T} + \frac{1}{N})\log(T)^2),$$

where $\mathbf{K}^*_{NT} = \text{diag}(1, \mathbf{K}_{NT})$ and $\boldsymbol{\beta}_i = (\rho_i, \mathbf{b}'_i)'$. Then,

$$\widetilde{\boldsymbol{\beta}}_i(\lambda_i) = \widetilde{\mathbf{X}}_i(\lambda_i)^{-1}\widetilde{\mathbf{X}}'_i\widetilde{\mathbf{X}}_i\mathbf{K}^*_{NT}\boldsymbol{\beta}_i + \widetilde{\mathbf{X}}_i(\lambda_i)^{-1}\widetilde{\mathbf{X}}'_i\begin{pmatrix} \mathbf{K}_{NT}\frac{1}{T}\sum_{t=1}^{T}\mathbf{f}_t\varepsilon_{i,t} \\ \mathbf{K}_{NT}\frac{1}{T}\sum_{t=2}^{T}\mathbf{f}_{t-1}\varepsilon_{i,t} \end{pmatrix} + \mathbf{R}_i, \qquad \text{(IA.21)}$$

10

where $\mathbf{R}_i$ is the remaining term, and we will show that $\sqrt{T}\mathbf{R}_i = o_p(1)$. Note that

$$\widetilde{\boldsymbol{\Sigma}}_{yf}(k) = \frac{1}{T} \sum_{t=k+1}^{T} \mathbf{y}_t \widehat{\mathbf{f}}'_{t-k} = \frac{1}{T} \sum_{t=k+1}^{T} \mathbf{y}_t \mathbf{f}'_{t-k} \mathbf{K}'_{NT} + \frac{1}{T} \sum_{t=k+1}^{T} \mathbf{y}_t (\widetilde{\mathbf{f}}_{t-k} - \mathbf{K}_{NT}\mathbf{f}_t)'$$
$$= \widehat{\boldsymbol{\Sigma}}_{yf} \mathbf{K}'_{NT} + O_p(N^{-1/2} + T^{-1/2})$$
$$\to_p \widehat{\boldsymbol{\Sigma}}_{yf} \mathbf{K}'_{NT}, \tag{IA.22}$$

if $N = o(T)$. Similarly, we can show that

$$\widetilde{\boldsymbol{\Sigma}}_f(k) = \mathbf{H}\boldsymbol{\Sigma}_f(k)\mathbf{H}' + o_p(1).$$

Therefore, if $N = o(T)$,

$$\widetilde{\mathbf{X}}_i = \begin{pmatrix} \widetilde{\boldsymbol{\Sigma}}'_{yf}\mathbf{w}_i & \widetilde{\boldsymbol{\Sigma}}'_f \\ \widetilde{\boldsymbol{\Sigma}}'_{yf}(1)\mathbf{w}_i & \widetilde{\boldsymbol{\Sigma}}'_f(1) \end{pmatrix} \to_p \mathbf{X}_i^H = \begin{pmatrix} \mathbf{H}\boldsymbol{\Sigma}'_{yf}\mathbf{w}_i & \mathbf{H}\boldsymbol{\Sigma}_f\mathbf{H}' \\ \mathbf{H}\boldsymbol{\Sigma}'_{yf}(1)\mathbf{w}_i & \mathbf{H}\boldsymbol{\Sigma}'_f(1)\mathbf{H}' \end{pmatrix},$$

and hence

$$\widehat{\mathbf{X}}'_i\widetilde{\mathbf{X}}_i \to_p \begin{pmatrix} \mathbf{w}'_i\boldsymbol{\Sigma}_{yf}\boldsymbol{\Sigma}'_{yf}\mathbf{w}_i + \mathbf{w}'_i\boldsymbol{\Sigma}_{yf}(1)\boldsymbol{\Sigma}'_{yf}(1)\mathbf{w}_i & \mathbf{w}'_i\boldsymbol{\Sigma}_{yf}\boldsymbol{\Sigma}_f\mathbf{H}' + \mathbf{w}'_i\boldsymbol{\Sigma}_{yf}(1)\boldsymbol{\Sigma}'_f(1)\mathbf{H}' \\ \mathbf{H}\boldsymbol{\Sigma}_f\boldsymbol{\Sigma}'_{yf}\mathbf{w}_i + \mathbf{H}\boldsymbol{\Sigma}_f(1)\boldsymbol{\Sigma}'_{yf}(1)\mathbf{w}_i & \mathbf{H}\boldsymbol{\Sigma}_f^2\mathbf{H}' + \mathbf{H}\boldsymbol{\Sigma}_f(1)\boldsymbol{\Sigma}'_f(1)\mathbf{H}' \end{pmatrix} = \mathbf{V}_i^H.$$

It follows that

$$(\widetilde{\mathbf{X}}'_i\widetilde{\mathbf{X}}_i + \lambda_i\mathbf{I}_{K+1})^{-1}\widetilde{\mathbf{X}}_i = O_p(1),$$

implying that $\sqrt{T}\mathbf{R}_i = o_p(1)$. Then (IA.21) implies that

$$\widetilde{\boldsymbol{\beta}}_i(\lambda_i) - \widetilde{\mathbf{X}}_i(\lambda_i)^{-1}\widetilde{\mathbf{X}}'_i\widetilde{\mathbf{X}}_i\mathbf{K}^*_{NT}\boldsymbol{\beta}_i = O_p(T^{-1/2}).$$

Let $\lambda_i \to 0$, we obtain that

$$\sqrt{T}(\widetilde{\mathbf{X}}'_i\widetilde{\mathbf{X}}_i)(\widetilde{\boldsymbol{\beta}}_i - \mathbf{K}^*_{NT}\boldsymbol{\beta}_i) = \widetilde{\mathbf{X}}'_i \begin{pmatrix} \mathbf{K}_{NT}\frac{1}{\sqrt{T}}\sum_{t=1}^{T} \mathbf{f}_t\varepsilon_{i,t} \\ \mathbf{K}_{NT}\frac{1}{\sqrt{T}}\sum_{t=2}^{T} \mathbf{f}_{t-1}\varepsilon_{i,t} \end{pmatrix} + o_p(1).$$

By a similar argument as that in the proof of Theorem 3, we have

$$\sqrt{T} \begin{pmatrix} \mathbf{K}_{NT}\frac{1}{T}\sum_{t=1}^{T} \mathbf{f}_t\varepsilon_{i,t} \\ \mathbf{K}_{NT}\frac{1}{T}\sum_{t=2}^{T} \mathbf{f}_{t-1}\varepsilon_{i,t} \end{pmatrix} \longrightarrow_d N(\mathbf{0}, \mathbf{U}_i^H), \tag{IA.23}$$

where

$$\mathrm{Var} \begin{pmatrix} \mathbf{K}_{NT}\frac{1}{\sqrt{T}}\sum_{t=1}^{T} \mathbf{f}_t\varepsilon_{i,t} \\ \mathbf{K}_{NT}\frac{1}{\sqrt{T}}\sum_{t=2}^{T} \mathbf{f}_{t-1}\varepsilon_{i,t} \end{pmatrix} \to \mathbf{U}_i^H,$$

which is defined as

$$\mathbf{U}_i^H = \begin{pmatrix} \mathbf{H}\boldsymbol{\Sigma}_{f\varepsilon_i}(0)\mathbf{H}' & \mathbf{H}\boldsymbol{\Sigma}_{f\varepsilon_i}(1)\mathbf{H}' \\ \mathbf{H}\boldsymbol{\Sigma}'_{f\varepsilon_i}(1)\mathbf{H}' & \mathbf{H}\boldsymbol{\Omega}_{f\varepsilon_i}(0)\mathbf{H}' \end{pmatrix},$$

where $\boldsymbol{\Sigma}_{f\varepsilon_i}(0)$, $\boldsymbol{\Sigma}_{f\varepsilon_i}(1)$, and $\boldsymbol{\Omega}_{f\varepsilon_i}(0)$ are defined in Section 4. It follows from (IA.23) that

$$\sqrt{T}\mathbf{V}_i^H(\widehat{\boldsymbol{\beta}}_i - \mathbf{K}_{NT}^*\boldsymbol{\beta}_i) \longrightarrow_d N(\mathbf{0}, \mathbf{X}_i^{H\prime}\mathbf{U}_i^H\mathbf{X}_i^H).$$

This completes the proof. $\square$