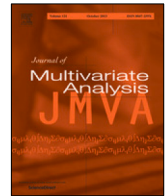




Contents lists available at ScienceDirect

Journal of Multivariate Analysis

journal homepage: www.elsevier.com/locate/jmva

A semiparametric latent factor model for large scale temporal data with heteroscedasticity

Lyuou Zhang^a, Wen Zhou^{b,*}, Haonan Wang^b

^a School of Statistics and Management, Shanghai University of Finance and Economics, Shanghai 200433, PR China

^b Department of Statistics, Colorado State University, Fort Collins, CO 80523, USA

ARTICLE INFO

Article history:

Received 22 December 2020

Received in revised form 5 July 2021

Accepted 6 July 2021

Available online 17 July 2021

AMS 2020 subject classifications:

primary 62H25

secondary 62M10

Keywords:

Efficient estimation

Heteroscedasticity

Large scale temporal data

Latent semiparametric factor model

Projection

ABSTRACT

Large scale temporal data have flourished in a vast array of applications, and their sophisticated structures, especially the heteroscedasticity among subjects with inter- and intra-temporal dependence, have fueled a great demand for new statistical models. In this paper, with covariate information, we consider a flexible model for large scale temporal data with subject-specific heteroscedasticity. Formally, the model employs latent semiparametric factors to simultaneously account for the subject-specific heteroscedasticity and the contemporaneous and/or serial correlations. The subject-specific heteroscedasticity is modeled as the product of the unobserved factor process and subject's covariate effect, which is further characterized via additive models. For estimation, we propose a two-step procedure. First, the latent factor process and nonparametric loading are recovered through projection-based methods, and following, we estimate the regression components by approaches motivated from the generalized least squares. By scrupulously examining the non-asymptotic rates for recovering the factor process and its loading, we show the consistency and efficiency of estimated regression coefficients in the absence of prior knowledge of latent factor process and subject's covariate effect. The statistical guarantees remain valid even for finite time points that makes our method particularly appealing when the subjects significantly outnumber the observation time points. Using comprehensive simulations, we demonstrate the finite sample performance of our method, which corroborates the theoretical findings. Finally, we apply our method to a data set of air quality and energy consumption collected at 129 monitoring sites in the United States in 2015.

© 2021 Elsevier Inc. All rights reserved.

1. Introduction

Jointly modeling a large and possibly divergent number of temporally evolving subjects arise ubiquitously in genomics, proteomics, environmental science, econometrics, clinical studies, and neuroscience. An extensively used statistical model for explaining the interactions and co-movements among the temporally evolving subjects is $y_{it} = \mathbf{z}_{it}^T \boldsymbol{\beta} + \varepsilon_{it}$, $i \in \{1, \dots, n\}$, $t \in \{1, \dots, T\}$, where y_{it} is the observation for the i th subject at time point t , $\boldsymbol{\beta}$ is a p -dimensional regression coefficient, \mathbf{z}_{it} is a p -dimensional covariate vector that might evolve in time, and $(\varepsilon_{1t}, \dots, \varepsilon_{nt})^T$ is a vector time series with possible contemporaneous correlations. Here, the number of subjects n is allowed to diverge much faster than the number of time points T . To name a few applications, y_{it} can model the expression level of the i th gene at time point t in

* Corresponding author.

E-mail addresses: zhanglvou@mail.shufe.edu.cn (L. Zhang), riczw@stat.colostate.edu (W. Zhou), wanghn@stat.colostate.edu (H. Wang).

a time course sequencing experiment, see, e.g., [10,40], the concentration of certain air pollutant in county i at day t , see, e.g., [29], and the measurement from electroencephalograms at brain location i and time point t in a motor-visual task experiment, see, e.g., [33]. As n rapidly grows, heteroscedasticity across subjects becomes inevitable and brings substantial challenges to modeling, estimation, and inference [14,46]. First, ignoring the subject-specific heteroscedasticity leads to inefficient estimation and inference on the regression coefficient β . In addition, in the presence of contemporaneous and serial correlations, when n rapidly outnumbers T the high dimensionality of data makes it even more difficult to accurately estimate the covariance of ϵ_{it} , which compromises the estimation efficiency of β .

In this paper, inspired by the approximate factor structure and its variants [11,26,38,44], we introduce a flexible data-driven model, where the heteroscedasticity across subjects and serial dependence of ϵ_{it} are assumed to arise from a product of the subject-specific effect and some latent stationary process. Specifically, motivated by Connor and Linton [13], Daye et al. [14], and Fan et al. [19], with additional time invariant covariates \mathbf{x}_i , we model the subject-specific effect by $\mathbf{g}(\mathbf{x}_i) = (g_1(\mathbf{x}_i), \dots, g_K(\mathbf{x}_i))^T$ with nonparametric functions g_1, \dots, g_K . For instance, \mathbf{x}_i could be the genetic information in a clinical study or the market capitalization in the financial asset allocation. Then, consider a K -dimensional zero-mean latent process \mathbf{f}_t , our semiparametric latent factor model with subject-specific heteroscedasticity is

$$y_{it} = \mathbf{z}_{it}^T \beta + \mathbf{g}(\mathbf{x}_i)^T \mathbf{f}_t + u_{it}, \quad (1)$$

where the residual process u_{it} is independent of \mathbf{f}_t . Analogous to the traditional factor models, $\mathbf{g}(\mathbf{x}_i)$ and \mathbf{f}_t can be viewed as the loading and factor, respectively. Particularly, $\mathbf{g}(\mathbf{x}_i)$ models the heteroscedasticity across subjects and, together with \mathbf{f}_t , retains the cross-sectional dependence while \mathbf{f}_t and u_{it} characterize the serial dependence. Model (1) features a large number of widely used models. For example, when \mathbf{f}_t is degenerate, (1) reduces to the partially linear additive models [41]; when $\mathbf{g}(\cdot)$ is known and \mathbf{f}_t is Gaussian, (1) is the traditional linear mixed model; with i replaced by one-dimensional spatial locations, (1) is analogous to the spatio-temporal model [31]; and when $\mathbf{g}(\cdot)$ reduces to constants, (1) is equivalent to the traditional factor models [2,11] or the panel data model with unobservable interactive effects [3,4]. It is worth mentioning that, focusing on the high-dimensional factor analysis, Connor and Linton [13] and Fan et al. [19] consider γ_{ik} in addition to $\mathbf{g}(\mathbf{x}_i)$ in the loading to account for that cannot be explained by \mathbf{x}_i . In contrast, one of our major goals is to efficiently estimate β in the presence of subject-specific heteroscedasticity and contemporaneous and serial correlations for any n and T . To that end, it requires recovering \mathbf{f}_t and $\mathbf{g}(\mathbf{x}_i)$ with satisfactory rates as well as accurate estimation of the long run variance of residues, while from Theorem 4.1 in Fan et al. [19] and Theorem 1, the desired rate for estimating the long run variance of residues will not be satisfied in the presence of γ_{ik} , unless more stringent conditions on n , T and γ_{ik} 's are imposed.

Like the partially linear model and the linear mixed model, regardless of its consistency, the ordinary least squares (OLS) estimator of β in (1) is inefficient without taking into account the dependence. Therefore, a careful estimation of the unobserved $\mathbf{g}(\mathbf{x}_i)$ and \mathbf{f}_t are needed to guarantee some sort of efficiency in both estimation and inference on β . In the literature, there are scatter approaches to estimate \mathbf{f}_t and its loadings for models similar to (1). For instance, Connor and Linton [13] employed a kernel method to estimate \mathbf{f}_t given \mathbf{x}_i with finite values, and Connor et al. [12] extended such estimates for general \mathbf{x}_i . Additionally, the consistency on estimated loading and latent factor shed light on estimating the large covariance matrix under assumptions of factor structures [18]. Motivated from these pioneering works, we propose a two-stage projection-based estimator of β , $\mathbf{g}(\mathbf{x}_i)$, and \mathbf{f}_t in (1). Roughly speaking, adopting a projection-based principal component type estimator [2,19], we first estimate $\mathbf{g}(\mathbf{x}_i)$ and \mathbf{f}_t from $y_{it} - \mathbf{z}_{it}^T \hat{\beta}^0$ for some preliminary $\hat{\beta}^0$. Using the estimated $\mathbf{g}(\mathbf{x}_i)$ and \mathbf{f}_t from the first-stage, we then estimate β with a generalized least squares (GLS) type approach.

Theoretically, the asymptotic properties such as consistency on estimating $\mathbf{g}(\mathbf{x}_i)$ and \mathbf{f}_t are not sufficient to guarantee the consistency and, especially the efficiency, of the second-stage estimator of β due to the lack of the finite sample characterization of errors from the first-stage [6]. In fact, it is known that a naive plug-in GLS type estimator does not necessarily guarantee the efficiency. To circumvent these challenges, a major contribution of this paper is a careful non-asymptotic analysis on the projection-based estimator of $\mathbf{g}(\mathbf{x}_i)$ and \mathbf{f}_t , by which we show that the consistency on estimating $\mathbf{g}(\mathbf{x}_i)$ and \mathbf{f}_t is free from restrictions on the relationship between n and T . With the exponential type concentration inequalities on estimating $\mathbf{g}(\mathbf{x}_i)$ and \mathbf{f}_t , we are able to obtain the finite sample deviation of the proposed two-stage estimator of β from the oracle GLS that enjoys full access to $\mathbf{g}(\mathbf{x}_i)$ and \mathbf{f}_t . These nontrivial results show that our estimator of β is overwhelmingly close to the oracle GLS, which establishes the efficiency of our estimator of β . In addition, we show the asymptotic normality of our estimator of β for drawing inference. The established concentration results for recovering $\mathbf{g}(\mathbf{x}_i)$ and \mathbf{f}_t (Theorem 1) are of independent interest for extending the projection-based principal component analysis (PCA) to other high-dimensional problems such as modeling temporally evolving tensor data or the segmentation of high-dimensional time series.

After introducing our model with identification conditions in Section 2.1, we detail the two-stage projection-based estimation of the loading, latent factor process, and regression coefficients in Section 2.2. We carry out the non-asymptotic analysis of our estimator and explore its efficiency in Section 3. Section 4 presents the inference. Sections 5 and 6 report simulation studies and an application on air quality and energy consumption data in the United States to demonstrate our method. After discussing potential extensions of our method in Section 7, we conclude the paper with technical details in Section 8. Extra numerical results and a detailed discussion on the determination of unknown dimension K are deferred to the supplement.

2. Methodology

2.1. A semiparametric latent factor model

Consider an $n \times 1$ vector of temporally evolving subjects $\mathbf{y}_t = (y_{1t}, \dots, y_{nt})^\top$ along with p -dimensional predictors \mathbf{z}_{it} and d -dimensional time invariant covariates \mathbf{x}_i associated with the i th subject. Our objective is to study the long run movement of y_{it} with respect to \mathbf{z}_{it} and model the dependence, over time, of each component of \mathbf{y}_t and across components, where the heteroscedasticity across subjects is accounted by \mathbf{x}_i . In our baseline formulation, each subject is modeled by a multi-factor linear model $y_{it} = \mathbf{z}_{it}^\top \boldsymbol{\beta} + \varepsilon_{it}$ [8] for $i \in \{1, \dots, n\}$ and $t \in \{1, \dots, T\}$, where $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$ is a p -dimensional vector of regression coefficients common across subjects. As discussed in Section 1, we adopt the semiparametric factor model $\varepsilon_{it} = \mathbf{g}(\mathbf{x}_i)^\top \mathbf{f}_t + u_{it}$, where the vector loading function $\mathbf{g}(\mathbf{x}_i) : \mathbb{R}^d \rightarrow \mathbb{R}^K$ accounts for the subject-specific heteroscedasticity and contemporaneous dependence and the K -dimensional latent factor process \mathbf{f}_t models the serial dependence. We further model $\mathbf{g}_k(\mathbf{x}_i)$ in an additive fashion, i.e., $\mathbf{g}_k(\mathbf{x}_i) = \sum_{\ell=1}^d g_{k\ell}(x_{i\ell})$, without losing flexibility yet providing concision in techniques [21]. Function \mathbf{g} offers structure that is flexible enough to allow dependence between $\{\mathbf{z}_{it}\}_{i \leq n, t \leq T}$ and $\{\mathbf{x}_i\}_{i \leq n}$, which has also been noted for other high-dimensional heteroscedastic regression models [14].

For each t , denote $\mathbf{Z}_t = (\mathbf{z}_{1t}, \dots, \mathbf{z}_{nt})^\top$, $\mathbf{u}_t = (u_{1t}, \dots, u_{nt})^\top$, and the $n \times K$ matrix of $\mathbf{g}_k(\mathbf{x}_i)$ by $\mathbf{G} = (\mathbf{g}(\mathbf{x}_1), \dots, \mathbf{g}(\mathbf{x}_n))^\top$. A more compact form of (1) reads

$$\mathbf{y}_t = \mathbf{Z}_t \boldsymbol{\beta} + \mathbf{G} \mathbf{f}_t + \mathbf{u}_t. \quad (2)$$

Similar to the traditional factor models, we impose the following conditions to control the rank and scale of latent loading function and factor process for model identification.

Condition 1. The rank of \mathbf{G} is K . For each t , $\mathbf{f}_t, \dots, \mathbf{f}_{Kt}$ are uncorrelated with each other and have zero mean and unit variance; u_{1t}, \dots, u_{nt} are uncorrelated with each other and have zero mean and finite variances; $\mathbf{Z}_t, \mathbf{f}_t$ and \mathbf{u}_t are uncorrelated from each other.

In addition, we assume that \mathbf{f}_t, u_{it} are independent from \mathbf{x}_i so that $\text{Cov}(y_{it}, y_{js} | \mathbf{x}_i, \mathbf{x}_j) = \mathbf{g}(\mathbf{x}_i)^\top \text{Cov}(\mathbf{f}_t, \mathbf{f}_s) \mathbf{g}(\mathbf{x}_j) + \text{Cov}(u_{it}, u_{js})$ for each i, j, t, s . Our model reaches beyond the existing literature [4,8] in the way that the inter- and intra-temporal dependence as well as the subject-specific heteroscedasticity are modeled simultaneously by \mathbf{f}_t and \mathbf{g} . Condition 1 is similar to conditions following (1.1) in Chamberlain and Rothschild [11] and Condition (C1) in Lam and Yao [26] (with diagonal $\boldsymbol{\Sigma}_\varepsilon$ and integer k herein). It guarantees the identifiability of the column space of \mathbf{G} . To further identify \mathbf{G} from its column space, consider T potentially dependent replicates $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_T)$ and $\mathbf{Z} = (\mathbf{Z}_1, \dots, \mathbf{Z}_T)$. Let $\mathbf{F} = (\mathbf{f}_1, \dots, \mathbf{f}_T)^\top$ and $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_T)$, (2) reads

$$\mathbf{Y} = \mathbf{Z}(\mathbf{I}_T \otimes \boldsymbol{\beta}) + \mathbf{G} \mathbf{F}^\top + \mathbf{U} \quad (3)$$

where \otimes is the Kronecker product. The following condition assures the identification of \mathbf{G} .

Condition 2. Almost surely, $T^{-1} \mathbf{F}^\top \mathbf{F} = \mathbf{I}_K$ and $\mathbf{G}^\top \mathbf{G}$ is diagonal with distinct entries.

Same as the PC1 condition of Bai and Ng [5], the first part of Condition 2 has been commonly assumed in factor analysis [20] and it is compatible with Condition 1 as $T^{-1} \mathbf{F}^\top \mathbf{F}$ is an estimator of $\text{Var}(\mathbf{f}_t)$. Under Condition 2, we can identify $\mathbf{G} \mathbf{H}$ and $\mathbf{F} \mathbf{H}$ for some $K \times K$ orthogonal matrix \mathbf{H} with $\mathbf{H} = \mathbf{I} + o(\min(n, T)^{-1})$, while the distinction among entries of $\mathbf{G}^\top \mathbf{G}$ further prevents rotational indeterminacy.

In contrast to the approximate factor model that allows cross-sectional dependence in \mathbf{u}_t , the assumption on u_{it} over i in Condition 1 is designated for efficiently estimating $\boldsymbol{\beta}$ without any restrictions on n and T . In fact, in the absence of the modulating component in (3), mild cross-sectional dependence of \mathbf{u}_t will not affect the estimation of \mathbf{G} and \mathbf{F} . On the other hand, without Condition 1 on u_{it} , a consistent estimate of $\text{Cov}(\mathbf{u}_t)$ is required for efficiently estimating $\boldsymbol{\beta}$. This will demand conditions on n and T , such as $\sqrt{n} \ln(n) = o(T)$ [18,43], which is more stringent compared to those in Section 3.

2.2. Two-stage projection-based estimation

2.2.1. First-stage: projection-based estimator of \mathbf{G} and \mathbf{F}

Let $\tilde{\mathbf{Y}} = \mathbf{Y} - \mathbf{Z}(\mathbf{I}_T \otimes \hat{\boldsymbol{\beta}}^0)$ and $\tilde{\mathbf{U}} = \mathbf{U} + \mathbf{Z}(\mathbf{I}_T \otimes (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^0))$ for some preliminary $\hat{\boldsymbol{\beta}}^0$, (3) can be written as

$$\tilde{\mathbf{Y}} = \mathbf{G} \mathbf{F}^\top + \tilde{\mathbf{U}}. \quad (4)$$

A naive approach is to estimate \mathbf{G} and \mathbf{F} using PCA. That is, \mathbf{F}/\sqrt{T} are estimated using eigenvectors corresponding to the first K largest eigenvalues of the $T \times T$ matrix $\tilde{\mathbf{Y}}^\top \tilde{\mathbf{Y}}$, and \mathbf{G} is estimated by right projecting $T^{-1} \tilde{\mathbf{Y}}$ onto the estimated \mathbf{F} . This method, however, takes into no account for the functional structure of \mathbf{g} or the smooth variation of $\{\tilde{y}_{it}\}_{i=1}^n$ from (4) against \mathbf{x}_i at each t . Fan et al. [19] proposed a projected principal component approach by smoothing $\{\tilde{y}_{it}\}_{i=1}^n$ as a function of \mathbf{x}_i at each t before implementing the aforementioned principal component estimation. Motivated by this, we replace $\tilde{\mathbf{Y}}$ by

$\tilde{\mathbf{P}}\tilde{\mathbf{Y}}$ for some projection \mathbf{P} onto a linear space spanned by a set of basis functions. Not only leveraging the smoothness, but \mathbf{P} can also be constructed to be orthogonal to errors $\tilde{\mathbf{U}}$ so that the subsequent PCA procedure is approximately errorless.

To begin with, let \mathbb{H} be a linear space spanned by a sequence of orthonormal basis functions $\{\phi_0(x) \equiv 1, \phi_1(x), \phi_2(x), \dots, \phi_J(x)\}$, where the number of basis functions J diverges in n . For each $k \in \{1, \dots, K\}$, $i \in \{1, \dots, n\}$, and $\ell \in \{1, \dots, d\}$, we have $g_{k\ell}(x_{i\ell}) = b_{0,k\ell} + \sum_{j=1}^J b_{j,k\ell} \phi_j(x_{i\ell}) + R_{k\ell}(x_{i\ell})$, where $\{b_{j,k\ell}\}_{j \leq J}$ are the coefficients and $R_{k\ell}$ is the approximation or projection error. Assume $Jd + 1 < n$ so that the coefficients are estimable. Denote $\mathbf{b}_k = (b_{0,k}, b_{1,k1}, \dots, b_{J,k1}, \dots, b_{1,kd}, \dots, b_{J,kd})^\top$ for each $k \in \{1, \dots, K\}$ and $i \in \{1, \dots, n\}$, where $b_{0,k} = \sqrt{J} \sum_{\ell=1}^d b_{0,k\ell}$, and $\boldsymbol{\varphi}_i = (1/\sqrt{J}, \phi_1(x_{i1}), \dots, \phi_J(x_{i1}), \dots, \phi_1(x_{id}), \dots, \phi_J(x_{id}))^\top$. Then, it admits $g_k(\mathbf{x}_i) = \boldsymbol{\varphi}_i^\top \mathbf{b}_k + \sum_{\ell=1}^d R_{k\ell}(x_{i\ell})$ and (4) can be rewritten as $\tilde{\mathbf{Y}} = (\tilde{\boldsymbol{\Phi}}\mathbf{B} + \mathbf{R})\mathbf{F}^\top + \tilde{\mathbf{U}}$, where $\tilde{\boldsymbol{\Phi}} = (\boldsymbol{\varphi}_1, \dots, \boldsymbol{\varphi}_n)^\top$, $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_K)$, and $\mathbf{R} = \{\sum_{\ell=1}^d R_{k\ell}(x_{i\ell})\}_{i=1, k=1}^{n,K}$. Then, we let $\mathbf{P} = \tilde{\boldsymbol{\Phi}}(\tilde{\boldsymbol{\Phi}}^\top \tilde{\boldsymbol{\Phi}})^{-1} \tilde{\boldsymbol{\Phi}}^\top$ and apply the PCA procedure to $\tilde{\mathbf{P}}\tilde{\mathbf{Y}}$. That is, we estimate $\hat{\mathbf{F}}$ by letting $\hat{\mathbf{F}}/\sqrt{T}$ be the eigenvectors corresponding to the first K largest eigenvalues of $\tilde{\mathbf{Y}}^\top \tilde{\mathbf{P}}\tilde{\mathbf{Y}}$ and estimate \mathbf{G} by $\hat{\mathbf{G}} = T^{-1} \tilde{\mathbf{P}}\tilde{\mathbf{Y}}\hat{\mathbf{F}}$. Moreover, we let $\hat{\mathbf{B}} = (\tilde{\boldsymbol{\Phi}}^\top \tilde{\boldsymbol{\Phi}})^{-1} \tilde{\boldsymbol{\Phi}}^\top \tilde{\mathbf{Y}}\hat{\mathbf{F}}$.

2.2.2. Second-stage: GLS-type estimator of $\boldsymbol{\beta}$

First, from (3), consider $\tilde{\mathbf{y}} = \mathbb{Z}_0^\top \boldsymbol{\beta} + \mathbf{G}^\top \sum_{t=1}^T \mathbf{f}_t + T^{-1} \sum_{t=1}^T \mathbf{u}_t$, where $\mathbb{Z}_0 = T^{-1} \sum_{t=1}^T \mathbf{Z}_t$ and $\tilde{\mathbf{y}} = T^{-1} \sum_{t=1}^T \mathbf{y}_t$. Conditional on \mathbf{Z}_t 's and \mathbf{x}_t 's, Condition 1 implies that the variance of $n \times 1$ vector $\tilde{\mathbf{y}}$ is

$$\mathbf{V} = \mathbf{G} \text{Var} \left(\frac{1}{T} \sum_{t=1}^T \mathbf{f}_t \right) \mathbf{G}^\top + \mathcal{D}, \quad (5)$$

where the $n \times n$ diagonal matrix \mathcal{D} has diagonal entries $\text{Var}(T^{-1} \sum_{t=1}^T u_{1t}), \dots, \text{Var}(T^{-1} \sum_{t=1}^T u_{nt})$. Then, (5) naturally leads to the oracle GLS-type estimate of $\boldsymbol{\beta}$,

$$\tilde{\boldsymbol{\beta}} = (\mathbb{Z}_0^\top \mathbf{V}^{-1} \mathbb{Z}_0)^{-1} \mathbb{Z}_0^\top \mathbf{V}^{-1} \tilde{\mathbf{y}}. \quad (6)$$

With the full knowledge on \mathbf{G} and \mathbf{F} in (3), \mathbf{V} in (6) can be estimated as follows. Let $\tilde{\mathbf{f}} = T^{-1} \sum_{t=1}^T \mathbf{f}_t$, it is known that $\text{Var}(T^{-1} \sum_{t=1}^T \mathbf{f}_t) = T^{-2} \sum_{t=-T+1}^{T-1} (T - |t|) \boldsymbol{\Sigma}_f(t)$, where $\boldsymbol{\Sigma}_f(s) = \text{Cov}(\mathbf{f}_t, \mathbf{f}_{t+s})$ and $\boldsymbol{\Sigma}_f(-s) = \text{Cov}(\mathbf{f}_{t-s}, \mathbf{f}_t)$ can be estimated by $\hat{\boldsymbol{\Sigma}}_f(s) = (T - s)^{-1} \sum_{t=1}^{T-s} (\mathbf{f}_t - \tilde{\mathbf{f}})(\mathbf{f}_{t+s} - \tilde{\mathbf{f}})^\top$ and $\hat{\boldsymbol{\Sigma}}_f(-s) = (T - s)^{-1} \sum_{t=s}^T (\mathbf{f}_{t-s} - \tilde{\mathbf{f}})(\mathbf{f}_t - \tilde{\mathbf{f}})^\top$ for $s \geq 0$, respectively. Naturally, $\nu(\mathbf{f}_t) = T^{-2} \sum_{t=-T+1}^{T-1} (T - |t|) \hat{\boldsymbol{\Sigma}}_f(t)$ estimates $\text{Var}(T^{-1} \sum_{t=1}^T \mathbf{f}_t)$ in (5), and similarly $\nu(u_{it})$ estimates $\text{Var}(T^{-1} \sum_{t=1}^T u_{it})$ for each $i \in \{1, \dots, n\}$ and the $n \times n$ diagonal matrix with diagonals $\nu(u_{1t}), \dots, \nu(u_{nt})$ estimates \mathcal{D} .

Though the oracle GLS estimator $\tilde{\boldsymbol{\beta}}$ is not accessible in practice as it depends on the full knowledge on \mathbf{f}_t and \mathbf{u}_t , it suggests a refined estimator of $\boldsymbol{\beta}$ by replacing \mathbf{G} and \mathbf{F} with $\hat{\mathbf{G}}$ and $\hat{\mathbf{F}}$ in (6), respectively. With $\hat{\mathbf{F}}$ from the first-stage, we can approximate $\nu(\mathbf{f}_t)$ and $\nu(u_{it})$ in (5) by $\nu(\hat{\mathbf{f}}_t)$ and $\nu(\hat{u}_{it})$ respectively, where $\hat{\mathbf{f}}_t$ is the t th row of $\hat{\mathbf{F}}$, and \hat{u}_t is the t th column of corresponding $\hat{\mathbf{U}} = \tilde{\mathbf{Y}} - \hat{\mathbf{G}}\hat{\mathbf{F}}^\top$. In summary, we can estimate \mathbf{V} by

$$\hat{\mathbf{V}} = \hat{\mathbf{G}}\nu(\hat{\mathbf{f}}_t)\hat{\mathbf{G}}^\top + \hat{\mathcal{D}}, \quad (7)$$

where $\hat{\mathcal{D}}$ is the $n \times n$ diagonal matrix with diagonals $\nu(\hat{u}_{1t}), \dots, \nu(\hat{u}_{nt})$ and we arrive at the **Two-stage Projection-based Estimator (TOPE)** of $\boldsymbol{\beta}$, $\hat{\boldsymbol{\beta}} = (\mathbb{Z}_0^\top \hat{\mathbf{V}}^{-1} \mathbb{Z}_0)^{-1} \mathbb{Z}_0^\top \hat{\mathbf{V}}^{-1} \tilde{\mathbf{y}}$. The procedure is summarized in Algorithm 1 below.

Algorithm 1. TOPE (Two-stage projection-based estimator)

Input: Data $\{(\mathbf{y}_{it}, \mathbf{x}_i, \mathbf{Z}_{it})\}_{i=1, t=1}^{n, T}$, predetermined K , and matrix of basis functions $\boldsymbol{\Phi}$.

Procedure:

1. For preliminary $\hat{\boldsymbol{\beta}}^0$, compute $\tilde{\mathbf{Y}} = \mathbf{Y} - \mathbb{Z}(\mathbf{I}_T \otimes \hat{\boldsymbol{\beta}}^0)$.
2. First-stage: estimate \mathbf{F} by letting the columns of $\hat{\mathbf{F}}/\sqrt{T}$ be the eigenvectors corresponding to the first K largest eigenvalues of $\tilde{\mathbf{Y}}^\top \tilde{\mathbf{P}}\tilde{\mathbf{Y}}$ and estimate \mathbf{G} by $\hat{\mathbf{G}} = \tilde{\mathbf{P}}\tilde{\mathbf{Y}}\hat{\mathbf{F}}$.
3. Second-stage: compute $\hat{\mathbf{V}} = \hat{\mathbf{G}}\nu(\hat{\mathbf{f}}_t)\hat{\mathbf{G}}^\top + \hat{\mathcal{D}}$ as in (7), where $\hat{\mathbf{f}}_t$ is the t th row of $\hat{\mathbf{F}}$ and \hat{u}_t is the t th column of $\hat{\mathbf{U}}$, and calculate **TOPE** $\hat{\boldsymbol{\beta}}$.

Output: $\hat{\mathbf{F}}, \hat{\mathbf{G}}, \hat{\boldsymbol{\beta}}$, and $\hat{\mathbf{V}}$.

The choice of preliminary $\hat{\boldsymbol{\beta}}^0$ for TOPE is quite flexible. Theoretically, the consistency and efficiency of TOPE are guaranteed whenever $\|\hat{\boldsymbol{\beta}}^0 - \boldsymbol{\beta}\|_2 = O_p(n^{-1/2+\alpha}T^{-1/2})$ for $\alpha \in [0, 1/2)$, which is not difficult to acquire. Some concrete choices of $\hat{\boldsymbol{\beta}}^0$ are discussed in depth in Section 8.3. Alternative to TOPE, one can first project \mathbf{Y} using $(\mathbf{I}_n - \mathbf{P})$, where $\mathbf{P} = \boldsymbol{\Phi}(\boldsymbol{\Phi}^\top \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top$. This leads to $(\mathbf{I}_n - \mathbf{P})\mathbf{Y} = (\mathbf{I}_n - \mathbf{P})\mathbb{Z}(\mathbf{I}_T \otimes \boldsymbol{\beta}) + \mathbf{R}\mathbf{F}^\top + (\mathbf{I}_n - \mathbf{P})\mathbf{U}$, which is similar to the procedure of profile likelihood [16] or restricted maximum likelihood [24]. However, the validity of this approach relies on the assumption that \mathbb{Z} and $\boldsymbol{\Phi}$ are linearly independent, which is more restricted than that of TOPE. Another seemingly straightforward approach is to project \mathbf{Y} using $(\mathbf{I}_n - \mathbf{P}_Z)$ where $\mathbf{P}_Z = \mathbb{Z}(\mathbb{Z}^\top \mathbb{Z})^{-1} \mathbb{Z}^\top$ and perform PCA on $(\mathbf{I}_n - \mathbf{P}_Z)\mathbf{Y} \approx (\mathbf{I}_n - \mathbf{P}_Z)\mathbf{G}\mathbf{F}^\top$ to estimate the loading and latent process. Though such an estimate of \mathbf{F} remains consistent, as noted by Wang et al. [45],

this approach only identifies the part of the latent structure that is orthogonal to \mathbb{Z} . That is, one can only obtain a consistent estimate of $(\mathbf{I}_n - \mathbf{P}_Z)\mathbf{G}$, and in particular $\widehat{\mathbf{G}} + \widetilde{\mathbf{P}}_Z\mathbf{A}$ will also be a valid estimator for arbitrary $n \times K$ matrix \mathbf{A} .

3. Theoretical properties of TOPE

3.1. Preliminaries

We impose the following conditions on our model, in addition to [Condition 1](#).

Condition 3. For any $\delta > 0$, $1 - n^{-1} \ln(1/\delta) \lesssim \lambda_{\min}(n^{-1}\mathbf{G}^\top \mathbf{G}) \leq \lambda_{\max}(n^{-1}\mathbf{G}^\top \mathbf{G}) \lesssim 1 + n^{-1} \ln(1/\delta)$ with probability at least $1 - \delta$.

Condition 4. The density of $\mathbf{x}_i \in \mathcal{X}^d$, where $\mathcal{X} \subset \mathbb{R}$ is compact, is bounded away from zero and infinity.

Condition 5 (Accuracy of the sieve approximation).

- (i) For each $\ell \in \{1, \dots, d\}$, $k \in \{1, \dots, K\}$, the loading function $g_{k\ell}(\cdot)$ belongs to a Hölder class $\mathcal{G} = \{g : |g^{(r)}(s) - g^{(r)}(t)| \leq L|s - t|^\gamma\}$ for some $L > 0$.
- (ii) For $\kappa = 2(r + \gamma) \geq 4$, $\sup_{\mathbf{x} \in \mathcal{X}} |g_{k\ell}(\mathbf{x}) - \sum_{j=1}^J b_{k,j\ell} \phi_j(\mathbf{x})|^2 \lesssim J^{-\kappa}$.
- (iii) It admits $\max_{k,j,\ell} b_{k,j\ell}^2 < \infty$.

[Condition 3](#) is similar to the pervasive condition on loading matrix in the traditional factor model [\[38\]](#). Since $\mathbf{G}\mathbf{G}^\top$ and $\mathbf{G}^\top \mathbf{G}$ have their first K largest eigenvalues in common, the K largest eigenvalues of $\mathbf{G}^\top \mathbf{G}$ also diverge in n . This condition ensures that \mathbf{x}_i has non-vanishing explanatory power on loading so that $\mathbf{G}^\top \mathbf{G}$ has spiked eigenvalues. [Condition 4](#) is standard in the literature of nonparametric and semiparametric statistics [\[22,23,39\]](#). The accuracy of sieve approximation in [Condition 5](#) can be obtained by common basis such as polynomial or B-splines [\[19,30\]](#).

Condition 6. For each $i \in \{1, \dots, n\}$, \mathbf{z}_{it} is weakly stationary. Almost surely, for each T we have,

- (i) eigenvalues of $n^{-1}\mathbb{Z}_0^\top \mathbb{Z}_0$ are bounded away from 0 and infinity;
- (ii) $\|\mathbf{P}_Z \mathbf{G}\|_{\mathbb{F}} = O(n^\alpha)$ for some $\alpha \in [0, 1/2)$, where \mathbf{P}_Z is the projection matrix on \mathbb{Z}_0 .

[Condition 6\(i\)](#) is similar to the standard condition on the design matrix in linear model that $\mathbb{Z}_0^\top \mathbb{Z}_0/n$ converges in n . Similar to conditions for semiparametric models in Robinson [\[35\]](#), (ii) guarantees identifications between the parametric and nonparametric parts in our model. Particularly, it allows consistent identification of the regression component without enforcing independence between \mathbf{z}_{it} and \mathbf{x}_i . For instance, in [Section 8.3](#), it is employed to show the existence of a legitimate preliminary estimator $\widehat{\beta}^0$.

At last, we impose some widely-used conditions [\[2,38\]](#) regarding the serial dependence and stationarity on $\{\mathbf{f}_t, \mathbf{u}_t\}$ as well as their tail behavior. Denote $\mathcal{F}_{-\infty}^0$ and \mathcal{F}_T^∞ the σ -algebra generated by $\{\mathbf{f}_t, \mathbf{u}_t\} : t \leq 0\}$ and $\{\mathbf{f}_t, \mathbf{u}_t\} : t \geq T\}$, and recall the α -mixing coefficient as $\alpha(T) = \sup_{A \in \mathcal{F}_{-\infty}^0, B \in \mathcal{F}_T^\infty} |\Pr(A) \Pr(B) - \Pr(A \cap B)|$.

Condition 7 (Serial dependence, stationarity, and tail behavior).

- (i) $\{\mathbf{u}_t, \mathbf{f}_t\}_{t \leq T}$ are strictly stationary with zero mean and finite long run variances.
- (ii) There exist $r_1, C_1 > 0$ such that for all $T > 0$, $\alpha(T) < \exp(-C_1 T^{r_1})$.
- (iii) There exist $r_2, r_3 > 1$ with $r_1^{-1} + r_2^{-1} + r_3^{-1} > 1$ and $b_1, b_2 > 0$ such that for each i, k, t and any $s > 0$, $\Pr(|u_{it}| > s) \leq \exp\{-(s/b_1)^{r_2}\}$ and $\Pr(|f_{tk}| > s) \leq \exp\{-(s/b_2)^{r_3}\}$.

3.2. Statistical guarantees

To establish the statistical guarantees of TOPE, we first perform a non-asymptotic analysis of $\widehat{\mathbf{F}}$ and $\widehat{\mathbf{G}}$, then obtain the deviation between $\widehat{\beta}$ and the GLS estimator $\widetilde{\beta}$ to study the efficiency of TOPE.

Theorem 1. Suppose that [Conditions 1, 2, and 3–7](#) hold. Assume $Jd + 1 < n$ and $J = o(n^{1/2-\alpha})$ for $\alpha \in [0, 1/2)$. With probability at least $1 - \delta$ for any $\delta > 0$, we have

$$\begin{aligned} \frac{1}{T} \|\widehat{\mathbf{F}} - \mathbf{F}\|_{\mathbb{F}}^2 &\lesssim \left(\frac{1}{n} + \frac{p^2}{n^{1-2\alpha}T} + \frac{1}{J^\kappa} \right) \ln(1/\delta), \\ \frac{1}{n} \|\widehat{\mathbf{G}} - \mathbf{G}\|_{\mathbb{F}}^2 &\lesssim \left(\frac{J}{n^2} + \frac{p^2 J}{n^{1-2\alpha}T} + \frac{p^4 J}{n^{2-4\alpha}T^2} + \frac{1}{J^{\kappa-1}} \right) \{\ln(1/\delta)\}^2, \\ \|\widehat{\mathbf{B}} - \mathbf{B}\|_{\mathbb{F}}^2 &\lesssim \left(\frac{J}{n^2} + \frac{p^2 J}{n^{1-2\alpha}T} + \frac{p^4 J}{n^{2-4\alpha}T^2} + \frac{1}{J^{\kappa-1}} \right) \{\ln(1/\delta)\}^2. \end{aligned}$$

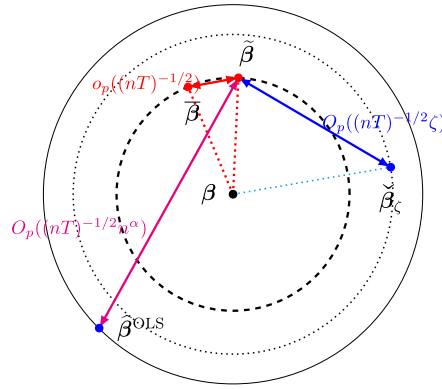


Fig. 1. A schematic illustration about different estimators to β in (1), where $\bar{\beta}$ is the TOPE estimator and $\tilde{\beta}$ is the oracle GLS estimator with full knowledge on \mathbf{G} and \mathbf{F} .

In contrast to the known asymptotic properties of $\hat{\mathbf{F}}$ and $\hat{\mathbf{G}}$ for traditional and semiparametric factor models with divergent n and T [5,19], [Theorem 1](#) provides finite sample characterization of \mathbf{F} and \mathbf{G} . Given a finite p , the rates obtained in [Theorem 1](#) agree with the asymptotic results in Fan et al. [19]. Also, whenever $p = o(n^{1/2-\alpha}T^{1/2}J^{-1/2})$, $\hat{\mathbf{F}}$ and $\hat{\mathbf{G}}$ are consistent in mean squared errors. Especially, for finite p , this consistency does not require T diverging to infinity which enables our method to be used for modeling a large number of short time series in practice. More importantly, results in [Theorem 1](#) make it possible to establish the following finite sample results on both β and its covariance with respect to the GLS estimator $\bar{\beta}$ as defined in (6).

Theorem 2. Under conditions in [Theorem 1](#), with probability at least $1 - \delta$,

$$\|\bar{\beta} - \tilde{\beta}\|_2 \lesssim \frac{1}{\sqrt{nT}} \left\{ \frac{\sqrt{J}}{n} + \frac{1}{\sqrt{n}} + \frac{1}{T} + \frac{p\sqrt{J}}{\sqrt{n^{1-2\alpha}T}} + \frac{1}{J^{(\kappa-1)/2}} \right\} \sqrt{\ln(1/\delta)},$$

where $\tilde{\beta}$ is the oracle GLS estimator of β with full knowledge of \mathbf{G} and \mathbf{F} as in [Section 2.2.2](#). In addition,

$$\|\text{Var}(\bar{\beta}) - \text{Var}(\tilde{\beta})\|_{\mathbb{F}} \lesssim \frac{p\vartheta_{n,T,J}}{nT} + \frac{p\vartheta_{n,T,J}^2}{(nT)^{3/2}},$$

where $\vartheta_{n,T,J} = n^{-1}J^{1/2} + n^{-1/2} + T^{-1} + pJ^{1/2}n^{-1/2+\alpha}T^{-1/2} + J^{-(\kappa-1)/2}$.

The nontrivial finite sample results in [Theorem 2](#) imply that the deviation between $\bar{\beta}$ and $\tilde{\beta}$ is due to: (i) the errors in estimating \mathbf{G} with rate $n^{-1/2}T^{-1/2}(n^{-1}J^{1/2} + pJ^{1/2}n^{-1/2+\alpha}T^{-1/2} + J^{-(\kappa-1)/2})$, (ii) the errors in estimating \mathbf{F} with rate $n^{-1/2}T^{-1/2}(n^{-1/2} + T^{-1} + pn^{-1/2+\alpha}T^{-1/2} + J^{-\kappa/2})$, and (iii) the deviation between $\text{Var}(\mathbf{f}_t)$ and $\text{Var}(T^{-1}\sum_{t=1}^T \mathbf{f}_t)$ with rate $n^{-1/2}T^{-3/2}$.

Let $\|\mathbf{A}\|_{\mathbf{S}} := n^{-1/2}\|\mathbf{S}^{-1/2}\mathbf{A}\mathbf{S}^{-1/2}\|_2$ for any positive definite \mathbf{S} and define the class of estimators to β with respect to working covariance \mathbf{V}_{ζ} by $\Theta_{\zeta} = \{\tilde{\beta}_{\zeta} = (\mathbf{Z}_0^{\top}\mathbf{V}_{\zeta}^{-1}\mathbf{Z}_0)^{-1}\mathbf{Z}_0^{\top}\mathbf{V}_{\zeta}^{-1}\tilde{\mathbf{y}} : \|\mathbf{V}_{\zeta} - \mathbf{V}\|_{\mathbf{V}} \lesssim \zeta\}$, where the oracle GLS estimator $\tilde{\beta} \in \Theta_0$, TOPE $\bar{\beta} \in \Theta_{\vartheta_{n,T,J}}$ by [Theorem 1](#), and OLS estimator $\tilde{\beta}^{\text{OLS}} \in \Theta_{n^{\alpha}}$ by [Proposition 5](#). From the proof of [Theorem 2](#), $\|\tilde{\beta}_{\zeta} - \tilde{\beta}\|_2 = O_p(n^{-1/2}T^{-1/2}\zeta)$ for any $\tilde{\beta}_{\zeta} \in \Theta_{\zeta}$. Thus, $\|\tilde{\beta}_{\zeta} - \tilde{\beta}\|_2 = O_p(n^{-1/2}T^{-1/2})$ if $\zeta = O(1)$ and $\|\tilde{\beta}_{\zeta} - \tilde{\beta}\|_2 = o_p(n^{-1/2}T^{-1/2})$ if $\zeta = o(1)$. In the presence of heteroscedasticity across subjects and/or autocorrelation, the oracle GLS estimator is known to be efficient in general [6]. Particularly, the GLS estimator $\tilde{\beta}$ is unbiased and efficient in Θ_{ζ} given the full information on \mathbf{G} and $\Sigma_f(t)$ for each $t \in \{1 - T, \dots, T - 1\}$. Therefore, [Theorem 2](#) implies that TOPE $\bar{\beta}$ is asymptotically unbiased, and given $p\vartheta_{n,T,J} = o(1)$, the non-asymptotic difference between the variances of $\bar{\beta}$ and $\tilde{\beta}$ is bounded by a rate smaller than $(nT)^{-1}$, which is the rate of $\text{Var}(\tilde{\beta})$. That is, the TOPE $\bar{\beta}$ is asymptotically efficient in Θ_{ζ} . This discussion is visualized in [Fig. 1](#).

Remark 1. Though [Theorems 1](#) and [2](#) are implicitly related to a legitimate preliminary estimator $\hat{\beta}^0$ via α , the existence of such a $\hat{\beta}^0$ is easily guaranteed under [Conditions 1, 2](#), and [6](#). In fact, as shown in [Proposition 5](#), the OLS estimator $\tilde{\beta}^{\text{OLS}}$ automatically satisfies $\|\tilde{\beta}^{\text{OLS}} - \beta\| = O_p(n^{-1/2+\alpha}T^{-1/2})$ and is therefore legitimate. Technically, $\hat{\beta}^0$ contributes to [Theorem 2](#) through $\vartheta_{n,T,J}$, more specifically via $p\sqrt{n^{-1/2+\alpha}T^{-1/2}}$. For a better $\hat{\beta}^0$ with decreasing α , its contribution through $\vartheta_{n,T,J}$ diminishes and achieves $p\sqrt{n^{-1/2}T^{-1/2}}$ when $\alpha \approx 0$. On the contrary, when $\alpha = 1/2 - 2c$ for any small $c > 0$, $\vartheta_{n,T,J} = o(n^{-c(\kappa-1)} + pn^{-c}T^{-1/2})$ and does not alter the conclusion in [Theorem 2](#). Hence, measured by α , $\hat{\beta}^0$ only affects on how $\vartheta_{n,T,J} = o(1)$ converges to zero as n and T increase, and it does not alter the conclusion about the efficiency of TOPE $\bar{\beta}$ as long as $\alpha \in [0, 1/2)$. These are also demonstrated by numerical studies in [Section A.4](#) in the supplement.

As a final remark, [Theorem 3](#) establishes results analogous to [Theorem 1](#) in the max norm and shares common observations with Wang and Fan [43] and Barigozzi et al. [7].

Theorem 3. For model (3), under the same conditions of [Theorem 1](#), with probability at least $1 - \delta$, $\|\widehat{\mathbf{F}} - \mathbf{F}\|_{\max} \lesssim \{n^{-1/2} + p^2(n^{1-\alpha}T)^{-1/2} + J^{-\kappa/2}\} \{\ln(T)\}^{2/3} \ln(1/\delta)$, $\|\widehat{\mathbf{G}} - \mathbf{G}\|_{\max} \lesssim \{[(T + p^2n^\alpha)\ln(n)T^{-2}]^{1/2} + n^{-1/2} + p^2(n^{1-\alpha}T)^{-1/2} + J^{-\kappa/2}\} \ln(1/\delta)$, and $\|\widehat{\mathbf{B}} - \mathbf{B}\|_{\max} \lesssim \{[(T + p^2n^\alpha)\ln(n)n^{-1}T^{-2}]^{1/2} + n^{-1/2} + p^2(n^{1-\alpha}T)^{-1/2} + J^{-\kappa/2}\} \ln(1/\delta)$.

4. TOPE-based inference

The following theorem provides inference on β based on TOPE. In general, the expectation with respect to $\{\mathbf{Z}_t, \mathbf{x}_i\}$ is unknown as their distribution is not accessible, therefore the asymptotic distribution of $\bar{\beta}$ conditional on $\{\mathbf{Z}_t, \mathbf{x}_i\}$ in (2) below is more practical.

Theorem 4. Under conditions in [Theorem 1](#), we have

- (i) with $\Sigma = E_{\mathbf{Z}_t, \mathbf{x}}\{(\mathbb{Z}_0^\top \mathbf{V}^{-1} \mathbb{Z}_0)^{-1}\}$, $\Sigma^{-1/2}(\bar{\beta} - \beta) \xrightarrow{d} N(\mathbf{0}, \mathbf{I}_p)$;
- (ii) conditional on \mathbf{Z}_t and \mathbf{x}_i , $(\mathbb{Z}_0^\top \mathbf{V}^{-1} \mathbb{Z}_0)^{1/2}(\bar{\beta} - \beta) \xrightarrow{d} N(\mathbf{0}, \mathbf{I}_p)$.

Replacing \mathbf{V} in [Theorem 4](#) (ii) by $\widehat{\mathbf{V}} = \widehat{\mathbf{G}}\widehat{\mathbf{V}}_t\widehat{\mathbf{G}}^\top + \widehat{\mathbf{D}}$ from (7), for any estimable $\mathbf{C}\beta$ with $q \times p$ matrix \mathbf{C} and $q < p$, $\text{CS}_{\mathbf{C}} = \{\mathbf{C}\beta : (\mathbf{C}\beta - \mathbf{C}\bar{\beta})^\top \{\mathbf{C}(\mathbb{Z}_0^\top \widehat{\mathbf{V}}^{-1} \mathbb{Z}_0)^{-1} \mathbf{C}^\top\}^{-1} (\mathbf{C}\beta - \mathbf{C}\bar{\beta}) < \chi_{q, 1-\eta}^2\}$ defines a $100(1-\eta)\%$ confidence set, where $\chi_{q, 1-\eta}^2$ is the $100(1-\eta)\%$ quantile of χ_q^2 distribution. When rows of \mathbf{C} are the natural basis of \mathbb{R}^p , $\text{CS}_{\mathbf{C}}$ provides a confidence set of a subset of β . Alternatively, denote $\widehat{\sigma}_\ell^2$ the ℓ th diagonal entry of $(\mathbb{Z}_0^\top \widehat{\mathbf{V}}^{-1} \mathbb{Z}_0)^{-1}$ for each $\ell \in \{1, \dots, p\}$, a $100(1-\eta)\%$ confidence interval for the ℓ th entry of β is $\text{CI}_\ell = [\bar{\beta}_\ell - \widehat{\sigma}_\ell \Phi^{-1}(1-\eta/2), \bar{\beta}_\ell + \widehat{\sigma}_\ell \Phi^{-1}(1-\eta/2)]$, where $\Phi(\cdot)$ is the cumulative distribution function of standard normal distribution. Moreover, [Theorem 4](#) implies that $\Pr(\|\bar{\beta} - \beta\|_\infty > \varepsilon) < p \exp(-\varepsilon^2 p^{-1} \sigma^{-2})$, where σ^2 can be estimated by the minimum diagonal of $(\mathbb{Z}_0^\top \widehat{\mathbf{V}}^{-1} \mathbb{Z}_0)^{-1}$. Thus, it also leads to a uniform confidence set for β at level $100(1-\eta)\%$, denoted as $\text{CI}' = \{\beta : |\beta_\ell - \bar{\beta}_\ell| \leq \widehat{\sigma} \sqrt{p \ln(p/\eta)}, \ell \in \{1, \dots, p\}\}$.

To draw inference on the explaining power of covariates \mathbf{x}_i on the dependence structure of data, Fan et al. [19] proposed a semiparametric specification testing statistic $S_G = \text{tr}\{(\widetilde{\mathbf{F}}^\top \widetilde{\mathbf{Y}}^\top \widetilde{\mathbf{Y}} \widetilde{\mathbf{F}})^{-1} \widetilde{\mathbf{F}}^\top \widetilde{\mathbf{Y}}^\top \widetilde{\mathbf{P}} \widetilde{\mathbf{Y}} \widetilde{\mathbf{F}}\}$, where $\widetilde{\mathbf{F}}/\sqrt{T}$ are the eigenvectors corresponding to the K largest eigenvalues of $\widetilde{\mathbf{Y}}^\top \widetilde{\mathbf{Y}}$. In addition to [Condition 3–7](#), assuming $T^{2/3} = o(n)$, $n\{\ln(n)\}^4 = o(T^2)$, $J = o(\min\{n^{1/2-\alpha}, \sqrt{T}\})$, and $\max\{T\sqrt{n}, n\} = o(J^\kappa)$, we have $(nS_G - JdK)(2JdK)^{-1/2} \xrightarrow{d} N(0, 1)$ whenever $\mathbf{G}(\mathbf{X}) = \mathbf{0}$. Thus, we can test $H_0 : \mathbf{G}(\mathbf{X}) = \mathbf{0}$ almost surely. Hence, S_G provides a diagnostic tool for the proposed model.

5. Numerical studies

5.1. Simulation settings

We demonstrate the finite sample performance of TOPE for both estimation and inference in comparison to three competing methods: the OLS estimator, which ignores heteroscedasticity across subjects and dependence; the GLS estimator, which naively utilizes the first K components of $T^{-1} \sum_{t=1}^T (\mathbf{y}_t - \bar{\mathbf{y}})(\mathbf{y}_t - \bar{\mathbf{y}})^\top$ as $\widehat{\mathbf{V}}$; and last, the oracle estimator, which is TOPE with known \mathbf{G} without using approximation. To implement TOPE, we employ the OLS estimator as preliminary $\widehat{\beta}^0$.

The mean squared error (MSE) and the empirical coverage probability (ECP) of the confidence region for β are employed to compare different methods. In addition, $\|\widehat{\mathbf{F}} - \mathbf{F}\|_F/\sqrt{T}$ and $\|\widehat{\mathbf{G}} - \mathbf{G}\|_F/\sqrt{n}$ are displayed to demonstrate estimations on \mathbf{G} and \mathbf{F} by TOPE. The maximum marginal length of the confidence set (MML) is used to demonstrate the efficiency. That is, the confidence set with ECP agreeing to the nominal level and small MML is preferable. For a clear presentation, we display MML of different methods normalized by the largest one (the MML of OLS, in general).

We consider $n \in \{50, 100, 200, 500, 1000, 2000\}$ and $T \in \{20, 50, 100, 200, 500\}$; also, we set $p = 4$ with $\beta = (1, 1, 1, 1)^\top$ and generate i.i.d. $z_{i\ell,t} \sim N(3\exp(t/30), 1)$ for each $i \in \{1, \dots, n\}$, $\ell \in \{1, \dots, p\}$, and $t \in \{1, \dots, T\}$. A similar setting was used in Huang et al. [23]. For the loading, we set $d = 3$ and generate i.i.d. $\mathbf{x}_i \sim U([0, 1]^d)$, then let $g_1(\mathbf{x}) = x_1$, $g_2(\mathbf{x}) = x_1^2 + x_2^2 - 1$, and $g_3(\mathbf{x}) = x_1^3 - 2x_1 + x_2$ for $K = 3$. As suggested by [19], with the initial realization \mathbf{G}_0 for g_1, g_2 and g_3 , we further compute $\mathbf{H}_G = \mathbf{G}_0^\top \mathbf{G}_0$ and set $\mathbf{G} = \mathbf{G}_0 \mathbf{H}_G$ in simulations so that [Condition 2](#) is satisfied.

The latent process $\mathbf{f}_t = (f_{1t}, f_{2t}, f_{3t})^\top$ consists of $K = 3$ independent univariate time series governed by the same model with one of the following three settings: independent in t , AR(1) with autoregressive coefficient $\rho = 0.5$, or ARMA(1, 1) with autoregressive coefficient $\rho = 0.5$ and moving average coefficient $\theta = 0.5$. Also, three innovations are considered: the standard normal, centered χ_5^2 , and t_8 . After generating \mathbf{f}_t , $\mathbf{F} = (\mathbf{f}_1, \dots, \mathbf{f}_T)^\top$ is further transformed so that $T^{-1} \mathbf{F}^\top \mathbf{F} = \mathbf{I}_K$ in [Condition 2](#) is satisfied. Similar to \mathbf{f}_t , we generate n independent \mathbf{u}_i from the same model, which includes two dependence structures: independent in t and AR(1) with autoregressive coefficient $\rho = 0.5$, as well as two innovations: $N(0, 0.01)$ and $(\chi_5^2 - 5)/10$. For each setting, 500 simulations are conducted.

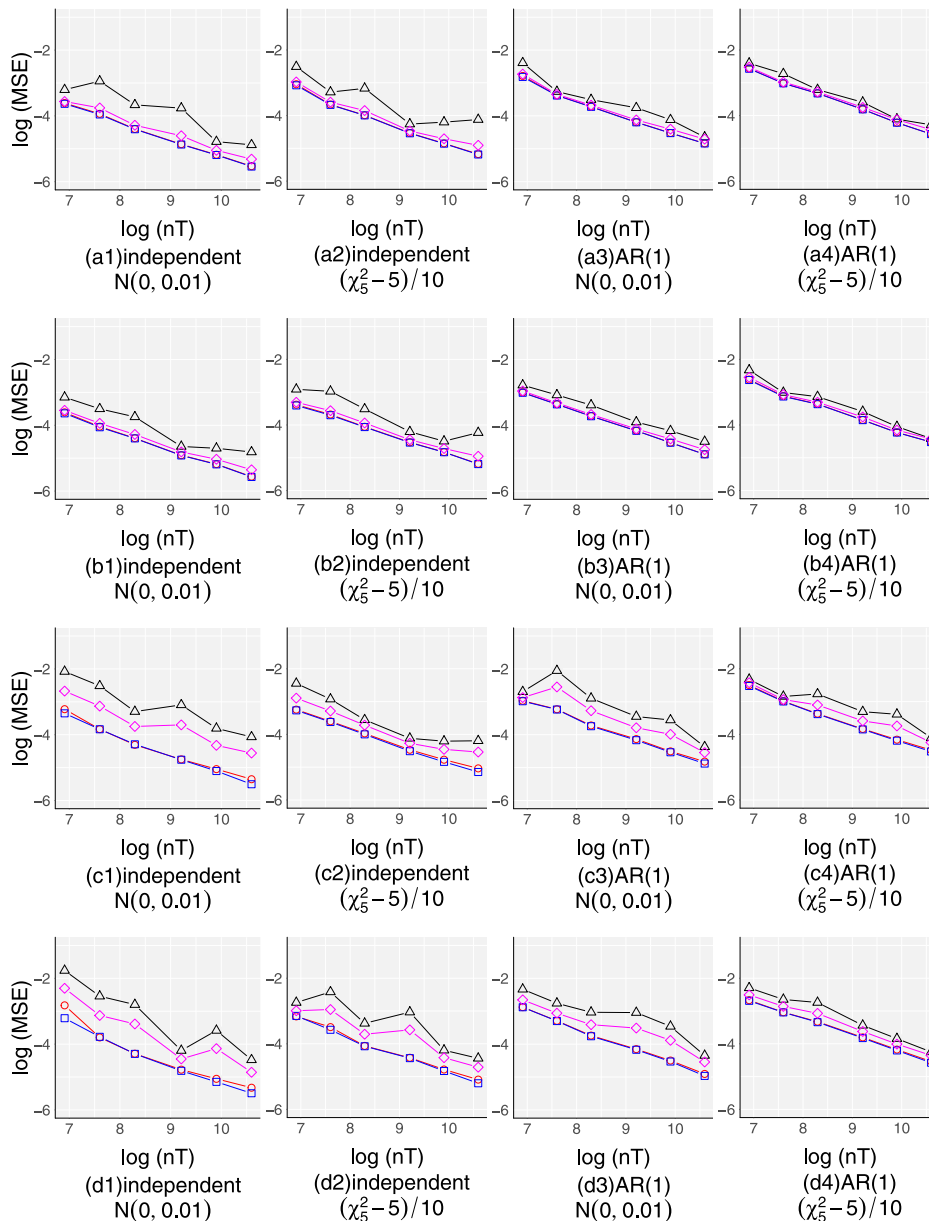


Fig. 2. Comparisons of the logarithm of MSE for estimating β by TOPE (“—○—”), along those of the oracle estimator (“—□—”), the GLS estimator (“—◇—”), and the OLS (“—△—”). Results are about $T = 20$. In plots (a1)–(a4), $f_{kt} \sim N(0, 1)$ is independent in k, t . In plots (b1)–(b4), $f_{kt} \sim t_8$ is independent in k, t . In plots (c1)–(c4), f_{kt} follows the ARMA(1, 1) model with $N(0, 1)$ innovation for each $k = 1, 2, 3$. In plots (d1)–(d4), f_{kt} follows the ARMA(1, 1) model with t_8 innovation for each $k = 1, 2, 3$. Distributions and serial correlations of \mathbf{u}_i are displayed in the plots.

5.2. Results

Fig. 2 displays the MSE with respect to $\ln(nT)$ on the logarithm scale when $T = 20$ and \mathbf{f}_t is independent in t or follows the ARMA(1, 1) model with $N(0, 1)$ or t_8 innovations. In Fig. 2, the MSEs of all estimators reduce as n increases. Both TOPE and GLS perform similarly as the oracle estimator when \mathbf{f}_t is independent in t ((a1)–(a4) and (b1)–(b4)), and outperform OLS; on the other hand, temporal dependence in \mathbf{u}_t slightly increases the MSE but does not alter the convergence rate ((c1)–(c4) and (d1)–(d4)). In the presence of temporal dependence in \mathbf{f}_t , GLS is outperformed while TOPE’s performance remains comparable to the oracle estimator ((c1)–(c4) and (d1)–(d4)). In the supplementary files, additional results for settings similar to Fig. 2 but with $T = 100, 500$ are reported in Figs. S.1–S.4, and results for \mathbf{f}_t following the AR(1) model

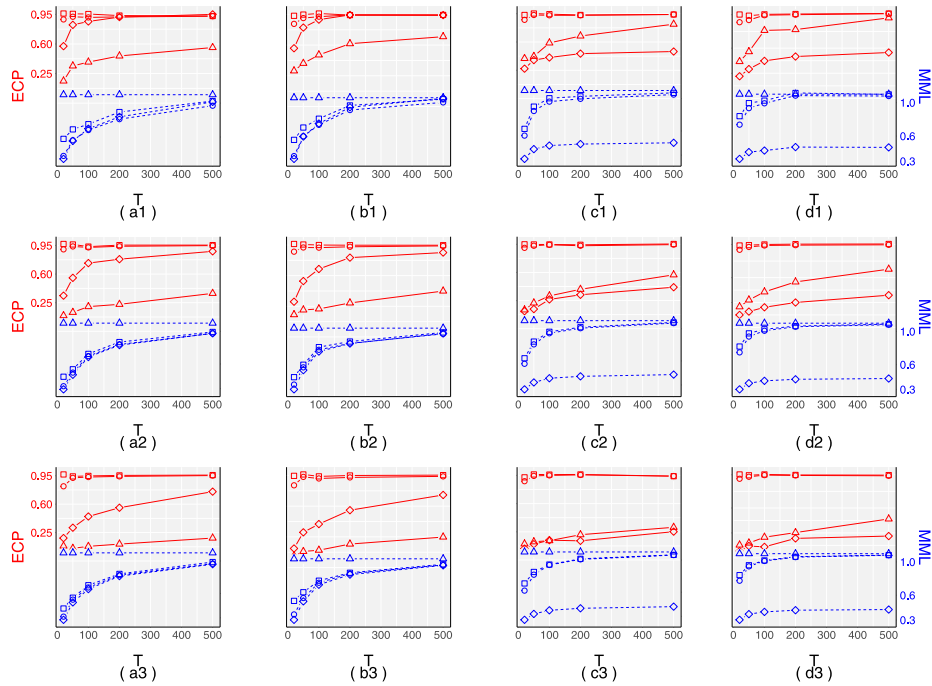


Fig. 3. Comparisons of the ECP and MML of 95% confidence region of TOPE (“—□—” for ECP and “—◇—” for MML) along those of the oracle estimator (“—□—” for ECP and “—◇—” for MML), the GLS estimator (“—◇—” for ECP and “—◇—” for MML), and the OLS (“—△—” for ECP and “—△—” for MML). In simulations, $f_{kt} \sim N(0, 1)$ is independent in k, t ; $n = 100, 500, 2000$ for the first, second, and third row, respectively. In plots (a1)–(a3) $u_{it} \sim N(0, 0.01)$ is independent in i, t . In plots (b1)–(b3), $u_{it} \sim (\chi^2_5 - 5)/10$ is independent in i, t . In plots (c1)–(c3) \mathbf{u}_i follows the AR(1) model with $N(0, 0.01)$ innovation while the same model is used for \mathbf{u}_i in plots (d1)–(d3) with $(\chi^2_5 - 5)/10$ innovation.

are displayed in Figs. S.5 and S.7. Similar observations are obtained for different settings of \mathbf{f}_t , and the differences among estimators decrease as T increases.

Fig. 3 displays the ECP and MML with respect to different T and n . The nominal level is 0.95. In Fig. 3, the confidence region of TOPE has ECP close to the nominal level with a small MML. Meanwhile, the coverage probabilities of OLS and GLS are both deviated from the nominal level and the deviations are substantial when n increases. In the presence of temporal dependence in \mathbf{u}_t , TOPE still outperforms GLS and OLS. The MML of TOPE substantially improves when n increases, particularly for large T , which reflects the fact that the estimation of \mathbf{F} in TOPE prefers large n (see (c1) and (c2), (d1)–(d2) in Fig. 3 for example). In the presence of the dependence of \mathbf{f}_t in t , TOPE performs remarkably well in terms of maintaining small MML and its ECP quickly converges to the nominal level in T (for example, Figs. S.17 and S.18 in the supplementary file). Meanwhile, due to the heteroscedasticity across subjects and the serial/cross-sectional correlations, both GLS and OLS fail to maintain the nominal coverage probability. As MML reflects the largest marginal variance of an estimator, OLS has large marginal variance in the presence of serial correlations in \mathbf{u}_t (Fig. 3 (d1)–(d3)). However, the ECP of OLS substantially deviates from the nominal level, which reflects the inconsistent covariance estimate of OLS. Also, it is interesting to notice that both the ECP and MML of GLS are smaller than those of others (Fig. 3 (c1)–(c3) and (d1)–(d3)), which shows that the naive GLS tends to ignore the serial correlations and greatly underestimate the variance that results in the poor confidence sets with low ECP. More simulation results are retained in the supplementary files and provide similar observations. Specifically, Figs. S.9 and S.12 in the supplementary files display results for f_{kt} independent in k, t with either independent u_{it} in i, t or \mathbf{u}_i following the AR(1) model with different innovations. Results for \mathbf{f}_t following the AR(1) model or the ARMA(1, 1) model with different innovations are included in Figs. S.13–S.24 in the supplementary files.

6. Study on air quality and energy consumption data using TOPE

In this section, we apply our method to analyze air quality data collected in the United States in 2015. The data consists of the mean PM2.5 concentration (in $\mu\text{g}/\text{m}^3$) from 129 monitoring sites on each Tuesday and Thursday in 2015 (<https://www.epa.gov/outdoor-air-quality-data>). We also include daily max 1-hour concentration of three common air pollutants, including NO_2 , SO_2 , and ozone, and the latitude and longitude of each monitoring site in our analysis. Sources of energy consumption are known as a potential factor to affect concentration of air pollutants. For this study, as covariates, we include the annual state-level energy consumption proportions of three major sources out of all possible resources,

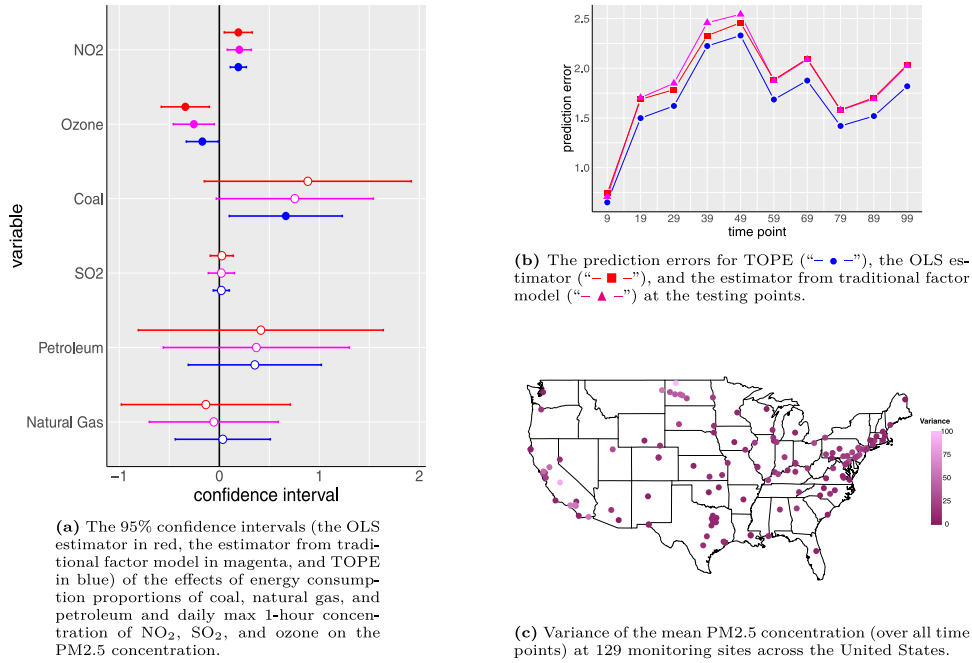


Fig. 4. Data display, resulting confidence intervals, and prediction comparison for the real data analysis. Panel (a), (b), and (c) are the 95% confidence intervals of the effects of different explanatory variables on the PM_{2.5} concentrations, the prediction errors for different methods, and the variance of the mean PM_{2.5} concentration across different monitoring sites, respectively.

namely coal, natural gas, and petroleum, in 2015 (<https://www.eia.gov/electricity/data/browser/>). For analysis, we log-transform the air pollutant data and remove potential seasonality. Also, we transform the latitude and longitude to keep their values within $[0, 1]$.

From Figs. 4(c) and S.29 in the supplement, it is observed that both geographical variables and energy consumption proportions help explaining the observed heteroscedasticity across monitoring sites so that we consider them as \mathbf{x}_i in (1). In this analysis, the daily max 1-hour concentration of NO₂, SO₂, and ozone, as well as the energy consumption proportions of coal, natural gas, and petroleum are considered as \mathbf{z}_{it} in (1).

To determine the dimension K of latent factor process, we apply both the eigenvalue-ratio procedure and the HDWN testing-based procedure proposed in the supplementary files. Ratios of the first ten adjacent eigenvalues of $\tilde{\mathbf{Y}}^T \tilde{\mathbf{P}} \tilde{\mathbf{Y}}$ are 4.13, 5.26, 6.58, 1.27, 1.68, 1.17, 1.29, 1.21, 1.10 such that the ratio between the third and fourth eigenvalues are the largest. On the other hand, for the HDWN testing-based procedure, the p -values for testing (B.1) with $K_0 = 1, 2$ and 3 are 0.026, 0.040 and 0.104, respectively. That is, we reject $H_0(1)$ and $H_0(2)$ but fail to reject $H_0(3)$ for (B.1). Thus, both eigenvalue-ratio procedure and the proposed HDWN testing-based procedure suggest $\bar{K} = 3$. Also, by the procedure discussed at the end in Section 4, we test $H_0 : \mathbf{G}(\mathbf{X}) = \mathbf{0}$ to further explore the statistical evidence to include geographical variables and energy consumption proportions to explain the heteroscedasticity across monitoring sites. We obtain $S_G = 2.34$ with p -value 4.84×10^{-14} ; thus, these covariates are included for modeling. Then, the complete model in the form of (1) for performing analysis on this data is

$$\ln(\text{PM}_{2.5, it}) = \beta_1 \ln(\text{NO}_{2, it}) + \beta_2 \ln(\text{SO}_{2, it}) + \beta_3 \ln(\text{Oz}_{it}) + \beta_4 \text{Cl}_i + \beta_5 \text{Ng}_i + \beta_6 \text{Pe}_i \\ + \sum_{k=1}^3 \{g_{k1}(\text{La}_i) + g_{k2}(\text{Lo}_i) + g_{k3}(\text{Cl}_i) + g_{k4}(\text{Ng}_i) + g_{k5}(\text{Pe}_i)\} f_{kt} + u_{it},$$

where $\ln(\text{PM}_{2.5, it})$ is the log concentration of PM_{2.5} from the monitoring site i at time t ; $\ln(\text{NO}_{2, it})$, $\ln(\text{SO}_{2, it})$, and $\ln(\text{Oz}_{it})$ are the log daily max 1-hour concentration of NO₂, SO₂, and ozone, respectively, from the same monitoring site i at time t ; Cl_i , Ng_i and Pe_i are the state-level energy consumption proportions of coal, natural gas, and petroleum out of all possible energy resources for the monitoring site i , respectively; and La_i and Lo_i are the latitude and longitude of the monitor site i , respectively.

For $g_{k\ell}$ in the above model, $\ell \in \{1, \dots, 5\}$, we use cubic spline with 11 knots to construct Φ for projection. We fit the above model using TOPE and draw inference as proposed in Section 4 to inspect the effects of covariates on the PM_{2.5} concentration. As an expected advantage, no further restrictions need to be imposed to model (1) and TOPE. In Fig. 4(a), the 95% confidence intervals for estimated coefficients using TOPE, the OLS estimator (by ignoring the variance

components), and the traditional factor model (using \mathbf{A} instead of $\mathbf{g}(\mathbf{x}_i)$ in the above model) are displayed for comparison. It reflects the efficiency of TOPE in the presence of heteroscedasticity across monitoring sites and serial/contemporaneous correlations discussed in Section 3. Specifically, the confidence intervals constructed by TOPE are the shortest among all three methods for all the six covariates. All methods suggest significant positive correlation between daily max 1-hour concentration of NO_2 and $\text{PM}_{2.5}$ concentration and significant negative correlation between ozone concentration and $\text{PM}_{2.5}$ concentration. TOPE reveals a significant positive correlation between coal consumption and $\text{PM}_{2.5}$ concentration, which agrees with Liang et al. [28] that coal consumption positively contributes to $\text{PM}_{2.5}$ concentration. However, this is missed by both the OLS estimator and the traditional factor model. In Fig. S.30 in the supplement, the recovered \mathbf{g}_k for $k \in \{1, 2, 3\}$ displays clear non-linearity.

In addition, we examine the prediction performance of TOPE, the OLS estimator, and the traditional factor model. For the 104 time points, we select 10 (from the 9th to the 99th, apart by 10 points) as the testing set and train the aforementioned three models using the remaining data points. With the estimated β from each method, the squared prediction errors at the testing points are displayed in Fig. 4(b). Compared with OLS and the traditional factor model, our model, alone with TOPE, has smaller prediction errors across all testing points, which demonstrates its superior prediction accuracy.

7. Discussions

Methodologically, we propose a flexible subject-specific heteroscedasticity model with latent semiparametric factor structures for analyzing large scale data with both intertemporal and intratemporal dependence. The model simultaneously accounts for the heteroscedasticity across subjects as well as the contemporaneous and serial correlations. We advocate a two-stage projection-based estimator for both the modulating and dependence components of the model, and establish an inference procedure for regression coefficients. We study the non-asymptotic rates for recovering the latent factor process and estimating the nonparametric loading function, which leads to the non-asymptotic properties of the estimated regression coefficients. As a result, we show that our proposed TOPE is asymptotically efficient within a fairly broad class of estimators including both the OLS and naive GLS estimators.

The widely-used Condition 2 essentially restricts \mathbf{F} to subspace $\{\mathbf{F} \in \mathbb{R}^{T \times K} : T^{-1}\mathbf{F}^\top \mathbf{F} = \mathbf{I}_K\}$, which might be stringent for some applications. In fact, we notice that it can be greatly relaxed by a concentration assumption of $T^{-1}\mathbf{F}^\top \mathbf{F}$ to \mathbf{I}_K , which can be derived from Condition 7 with the help of the so-called τ -mixing coefficient. As a result, this will alter the convergence rate of $\hat{\mathbf{F}} - \mathbf{F}$ in Theorem 1. Furthermore, as noted after Condition 2, we assume that the residual process u_{it} is uncorrelated over i to establish the statistical guarantee of TOPE on estimating β . This condition is similar to that of the traditional PCA that assumes uncorrelated samples. It can be further relaxed to, for example, $\max_{i \leq n} \sum_{t=1}^T |E(u_{it}u_{jt})| < C_2$, $\max_{i \leq n} \sum_{k=1}^n \sum_{m=1}^n \sum_{t=1}^T \sum_{s=1}^T |\text{cov}(u_{it}u_{kt}, u_{is}u_{ms})| < C_2$, and $(nT)^{-1} \sum_{i=1}^n \sum_{j=1}^n \sum_{t=1}^T \sum_{s=1}^T |E(u_{it}u_{js})| < C_2$ for some $C_2 > 0$. However, as a result, the $n \times n$ covariance matrix $\text{Cov}(\mathbf{u}_t)$ must be used in place of \mathcal{D} in (5) to retain the efficiency of TOPE. For that purpose, both the weighted PCA [25] and the estimator using thresholding principal orthogonal complements [18] can be employed in conjunction with TOPE. Then, in addition to some more stringent conditions on n and T , the non-asymptotic results must be re-established to obtain the similar conclusions in Section 3. Finally, from its construction, TOPE also paves a potentially effective way to model high-dimensional temporal data with multiple responses and simultaneously draw inference on the heteroscedasticity. We will explore these questions in future efforts.

8. Proofs of main theorems and technical results

We begin by presenting some notation. For a matrix $\mathbf{M} = (m_{ij})_{1 \leq i, j \leq p} \in \mathbb{R}^{p \times p}$, denote $\|\mathbf{M}\|_F = (\sum_{i=1}^p \sum_{j=1}^p m_{ij}^2)^{1/2}$ the Frobenius norm, $\|\mathbf{M}\|_{\max} = \max_{i,j} |m_{ij}|$ the maximum norm, and $\|\mathbf{M}\|_\infty = \max_i \sum_j |m_{ij}|$ the induced ℓ_∞ norm. The spectral norm of \mathbf{M} corresponds to its largest singular value, defined as $\|\mathbf{M}\|_2 = \sup_{\mathbf{a} \in S} \|\mathbf{M}\mathbf{a}\|_2$, where $S = \{\mathbf{a} \in \mathbb{R}^p : \|\mathbf{a}\|_2 = 1\}$ and the ℓ_q -norm of p -dimensional vector $\mathbf{a} = (a_1, \dots, a_p)^\top$ is defined by $\|\mathbf{a}\|_q = (\sum_{j=1}^p |a_j|^q)^{1/q}$ with $1 \leq q < \infty$. Denote the minimum and maximum eigenvalues of \mathbf{M} by $\lambda_{\min}(\mathbf{M})$ and $\lambda_{\max}(\mathbf{M})$, respectively. Let $\text{tr}(\mathbf{M}) = \sum_{j=1}^p m_{jj}$ and $\text{vec}(\mathbf{M})$ be the trace and vectorization of \mathbf{M} , and \otimes be the Kronecker product. We write \mathbf{I} for an identity matrix. For sequences $\{a_n\}$ and $\{b_n\}$, $a_n = o(b_n)$ if $a_n/b_n \rightarrow 0$ as $n \rightarrow \infty$ and $a_n = O(b_n)$ if $\limsup_{n \rightarrow \infty} |a_n|/b_n < \infty$; $X_n = o_p(a_n)$ and $X_n = O_p(a_n)$ are similarly defined for a sequence of random variables X_n ; $a_n \lesssim b_n$ if and only if $a_n \leq Cb_n$ for some C independent of n ; and $a_n \asymp b_n$ if and only if there exist C, D independent of n such that $C|b_n| \leq |a_n| \leq D|b_n|$. Denote \xrightarrow{p} and \xrightarrow{d} the convergence in probability and in distribution, respectively. Unless specified otherwise, $\delta > 0$ and $C > 0$ denote absolute constants independent of n, T, p .

Remark 2. The techniques in this section primarily depend on the derivation of a series of nontrivial exponential type concentration inequalities for preliminary estimators (such as $\hat{\beta}^0$ or $\hat{\mathbf{F}}$) and their approximations (such as $\hat{\mathbf{F}} - \mathbf{F}\mathbf{H} = \sum_{i=1}^8 (\mathbf{A}_i \mathbf{K}^{-1})$ before Lemma 6). Together with the union bounds, it avoids entangling with the correlations between any preliminary estimators and the data.

8.1. Proof of the main results

8.1.1. Invertibility of the projection matrix

Without loss of generality, we take $\mathcal{X}^d = [0, 1]^d$. Consider coefficients $\mathbf{a}_k = (a_0^{(k)}, a_{11}^{(k)}, \dots, a_{j_1}^{(k)}, \dots, a_{1d}^{(k)}, \dots, a_{jd}^{(k)})^\top \in \mathbb{R}^{jd+1}$ for $k \geq 1$, and define

$$\langle \mathbf{a}_1, \mathbf{a}_2 \rangle_n = \frac{1}{n} \sum_i \{a_0^{(1)} + \sum_j \sum_\ell a_{j\ell}^{(1)} \phi_j(X_{i\ell})\} \{a_0^{(2)} + \sum_j \sum_\ell a_{j\ell}^{(2)} \phi_j(X_{i\ell})\}. \quad (8)$$

In the literature, conditions on the largest and smallest eigenvalues of $n^{-1} \Phi^\top \Phi$ are usually stated as key assumptions for theoretical guarantees, see, e.g., [19]. Under standard nonparametric settings, we can establish it as follows.

Lemma 1. Under Condition 4, whenever $J = o(\sqrt{n})$ and $d < J^{-1}n$, with probability at least $1 - \delta$,

$$n \left\{ 1 - \frac{J}{n} \ln(J^2/\delta) \right\} \lesssim \lambda_{\min}(\Phi^\top \Phi) < \lambda_{\max}(\Phi^\top \Phi) \lesssim n \left\{ 1 + \frac{J}{n} \ln(J^2/\delta) \right\},$$

where, as defined in Section 2.2,

$$\Phi = \begin{bmatrix} 1/\sqrt{J} & \phi_1(x_{11}) & \dots & \phi_J(x_{11}) & \dots & \phi_1(x_{1d}) & \dots & \phi_J(x_{1d}) \\ \vdots & \vdots & & \vdots & & \vdots & & \vdots \\ 1/\sqrt{J} & \phi_1(x_{n1}) & \dots & \phi_J(x_{n1}) & \dots & \phi_1(x_{nd}) & \dots & \phi_J(x_{nd}) \end{bmatrix}.$$

Proof. From (8), $\langle \mathbf{a}, \mathbf{a} \rangle_n = \mathbf{a}^\top (n^{-1} \Phi^\top \Phi) \mathbf{a}$ for any $\mathbf{a} \in \mathbb{R}^{jd+1}$. For any $\delta > 0$, let $\mathcal{A}_\delta = \{|\langle \mathbf{a}, \mathbf{a} \rangle_n - E(\langle \mathbf{a}, \mathbf{a} \rangle_n)| \gtrsim n^{-1} J \ln(J^2/\delta) E(\langle \mathbf{a}, \mathbf{a} \rangle_n)\}$. On \mathcal{A}_δ^c , we have $\{1 - n^{-1} J \ln(J^2/\delta)\} E(\langle \mathbf{a}, \mathbf{a} \rangle_n) \lesssim \langle \mathbf{a}, \mathbf{a} \rangle_n \lesssim \{1 + n^{-1} J \ln(J^2/\delta)\} E(\langle \mathbf{a}, \mathbf{a} \rangle_n)$. By Lemma 2, $E(\langle \mathbf{a}, \mathbf{a} \rangle_n) \asymp \|\mathbf{a}\|_2^2$. Thus, $\{1 - n^{-1} J \ln(J^2/\delta)\} \|\mathbf{a}\|_2^2 \lesssim \mathbf{a}^\top (n^{-1} \Phi^\top \Phi) \mathbf{a} \lesssim \{1 + n^{-1} J \ln(J^2/\delta)\} \|\mathbf{a}\|_2^2$. The conclusion follows Lemma 3, which implies $\Pr\{\mathcal{A}_\delta\} < \delta$. \square

8.1.2. Proof of main theorems

Proof of Theorems 1 and 3. Theorems 1 and 3 readily follow from Propositions 1–4. \square

Proof of Theorem 2. Recall that $\widehat{\mathbf{V}} = \widehat{\mathbf{G}} \nu(\widehat{\mathbf{f}}_t) \widehat{\mathbf{G}}^\top + \widehat{\mathbf{D}}$, similarly to the proof of Lemma 14, we have

$$\lambda_{\min}(\widehat{\mathbf{V}}) \gtrsim \frac{1}{T} \left[1 + \frac{1}{\sqrt{nT}} + \sqrt{\frac{T + p^2 n^{2\alpha}}{n^3 T}} + \frac{\{(T + p^2 n^{2\alpha}) \ln(n)\}^{1/4}}{\sqrt{n^2 T}} + \frac{1}{nJ^{\kappa/2}} \right] \{1 + \sqrt{\ln(1/\delta)}\},$$

with probability at least $1 - \delta$. Then, by Lemma 14, with probability at least $1 - \delta$,

$$\begin{aligned} \|\widehat{\mathbf{V}}^{-1} - \mathbf{V}^{-1}\|_2 &= \|\widehat{\mathbf{V}}^{-1}(\mathbf{V} - \widehat{\mathbf{V}})\mathbf{V}^{-1}\|_2 \leq \|\widehat{\mathbf{V}}^{-1}\|_2 \|\mathbf{V}^{-1}(\widehat{\mathbf{V}} - \mathbf{V})\|_2 \\ &\lesssim T \left\{ \frac{\sqrt{J}}{n} + \frac{1}{\sqrt{n}} + \frac{1}{T} + \frac{p\sqrt{J}}{\sqrt{n^{1-2\alpha}T}} + \frac{1}{J^{(\kappa-1)/2}} \right\} \{1 + \sqrt{\ln(1/\delta)}\}. \end{aligned}$$

By Lemmas 16 and 17, $\|\mathbb{Z}_0^\top (\mathbf{G}\bar{\mathbf{f}} + \bar{\mathbf{u}})\|_2 \lesssim \|\mathbb{Z}_0\|_{\mathbb{F}} T^{-1/2} \sqrt{\ln(1/\delta)}$ with probability at least $1 - \delta$. Thus, with probability at least $1 - \delta$,

$$\begin{aligned} \|\bar{\boldsymbol{\beta}} - \widetilde{\boldsymbol{\beta}}\|_2 &\leq \|(\mathbb{Z}_0^\top \widehat{\mathbf{V}}^{-1} \mathbb{Z}_0)^{-1} - (\mathbb{Z}_0^\top \mathbf{V}^{-1} \mathbb{Z}_0)^{-1}\| \mathbb{Z}_0^\top \mathbf{V}^{-1} (\mathbf{G}\bar{\mathbf{f}} + \bar{\mathbf{u}})\|_2 \\ &\quad + \|(\mathbb{Z}_0^\top \widehat{\mathbf{V}}^{-1} \mathbb{Z}_0)^{-1} \mathbb{Z}_0^\top (\widehat{\mathbf{V}}^{-1} - \mathbf{V}^{-1}) (\mathbf{G}\bar{\mathbf{f}} + \bar{\mathbf{u}})\|_2 \\ &\lesssim \frac{1}{\sqrt{nT}} \left\{ \frac{\sqrt{J}}{n} + \frac{1}{\sqrt{n}} + \frac{1}{T} + \frac{p\sqrt{J}}{\sqrt{n^{1-2\alpha}T}} + \frac{1}{J^{(\kappa-1)/2}} \right\} \{1 + \sqrt{\ln(1/\delta)}\}. \end{aligned}$$

Therefore, for any $a > 0$,

$$\begin{aligned} E(\|\bar{\boldsymbol{\beta}} - \widetilde{\boldsymbol{\beta}}\|_2^2) &= \int_0^\infty \Pr(\|\bar{\boldsymbol{\beta}} - \widetilde{\boldsymbol{\beta}}\|_2^2 > s) ds = \int_0^a \Pr(\|\bar{\boldsymbol{\beta}} - \widetilde{\boldsymbol{\beta}}\|_2^2 > s) ds + \int_a^\infty \Pr(\|\bar{\boldsymbol{\beta}} - \widetilde{\boldsymbol{\beta}}\|_2^2 > s) ds \\ &\leq a + e \int_a^\infty \exp\{-snT(C\vartheta_{n,T,J}^2)^{-1}\} ds \\ &= a + Ce\vartheta_{n,T,J}^2(nT)^{-1} \exp\{-anT(C\vartheta_{n,T,J}^2)^{-1}\}, \end{aligned}$$

with $\vartheta_{n,T,J} = J^{1/2}n^{-1} + n^{-1/2} + T^{-1} + pJ^{1/2}n^{-1/2+\alpha}T^{-1/2} + J^{-(\kappa-1)/2}$ and constant $C > 0$. Letting $a = (nT)^{-1}C\vartheta_{n,T,J}^2$ gives $E(\|\bar{\boldsymbol{\beta}} - \widetilde{\boldsymbol{\beta}}\|_2^2) \leq 2(nT)^{-1}C\vartheta_{n,T,J}^2$. For TOPE $\bar{\boldsymbol{\beta}}$ and the oracle GLS estimator $\widetilde{\boldsymbol{\beta}}$ whose j th components are denoted by $\bar{\beta}_j$ and

$\tilde{\beta}_j$ respectively, repeatedly employing Cauchy–Schwarz inequality to each of the (i, j) pair with $i, j \in \{1, \dots, p\}$ leads to

$$\begin{aligned} |\text{Cov}(\tilde{\beta}_i, \tilde{\beta}_j) - \text{Cov}(\tilde{\beta}_i, \tilde{\beta}_j)| &= |\mathbb{E}\{(\tilde{\beta}_i - \beta_i)(\tilde{\beta}_j - \beta_j)\} + \mathbb{E}\{(\tilde{\beta}_i - \tilde{\beta}_i)(\tilde{\beta}_j - \beta_j)\}| \\ &\leq [\mathbb{E}\{(\tilde{\beta}_i - \beta_i)^2\}]^{1/2} [\mathbb{E}\{(\tilde{\beta}_j - \beta_j)^2\}]^{1/2} + [\mathbb{E}\{(\tilde{\beta}_j - \beta_j)^2\}]^{1/2} [\mathbb{E}\{(\tilde{\beta}_i - \tilde{\beta}_i)^2\}]^{1/2} \\ &\lesssim \frac{\vartheta_{n,T,J}}{nT} + \frac{\vartheta_{n,T,J}^2}{(nT)^{3/2}}, \end{aligned}$$

which yields

$$\|\text{Var}(\tilde{\beta}) - \text{Var}(\tilde{\beta})\|_{\mathbb{F}} = \left[\sum_{i,j=1}^p \{\text{Cov}(\tilde{\beta}_i, \tilde{\beta}_j) - \text{Cov}(\tilde{\beta}_i, \tilde{\beta}_j)\}^2 \right]^{1/2} \lesssim \frac{p\vartheta_{n,T,J}}{nT} + \frac{p\vartheta_{n,T,J}^2}{(nT)^{3/2}}. \quad \square$$

Proof of Theorem 4. (i) For the oracle GLS estimator $\tilde{\beta}$, it holds

$$\tilde{\beta} - \beta = (\mathbb{Z}_0^\top \mathbf{V}^{-1} \mathbb{Z}_0)^{-1} \mathbb{Z}_0^\top \mathbf{V}^{-1} \left(\mathbf{G} \frac{1}{T} \sum_{t=1}^T \mathbf{f}_t + \frac{1}{T} \sum_{t=1}^T \mathbf{u}_t \right) := \mathbf{A} \left(\mathbf{G} \frac{1}{T} \sum_{t=1}^T \mathbf{f}_t + \frac{1}{T} \sum_{t=1}^T \mathbf{u}_t \right),$$

where $\mathbf{A} = (\mathbb{Z}_0^\top \mathbf{V}^{-1} \mathbb{Z}_0)^{-1} \mathbb{Z}_0^\top \mathbf{V}^{-1}$ and $\mathbf{V} = \mathbf{G} \text{Var}(T^{-1} \sum_{t=1}^T \mathbf{f}_t) \mathbf{G}^\top + \text{Var}(T^{-1} \sum_{t=1}^T \mathbf{u}_t) \mathbf{I}_n$ as defined in (5). For any p -vector \mathbf{c} , $(nT)^{1/2} \mathbf{c}^\top (\tilde{\beta} - \beta) := \sum_{t=1}^T (W_{nt} + \tilde{W}_{nt})$, where $W_{nt} = n^{1/2} T^{-1/2} \mathbf{c}^\top \mathbf{A} \mathbf{u}_t$ and $\tilde{W}_{nt} = n^{1/2} T^{-1/2} \mathbf{c}^\top \mathbf{A} \mathbf{G} \mathbf{f}_t$. Then $\sum_{t=1}^T \mathbb{E}[|W_{nt}|^3] = n^{3/2} T^{-1/2} \|\mathbf{c}\|_2^3 \mathbb{E}[\|\mathbf{A}\|_2^3] \mathbb{E}[\|\mathbf{u}_1\|_2^3] < \infty$ for any n , and since $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_{\mathbb{F}} \leq p \|\mathbf{A}\|_2$ we have

$$\frac{\sum_{t=1}^T \mathbb{E}[|W_{nt}|^3]}{(\sum_{t=1}^T \mathbb{E}[W_{nt}^2])^{3/2}} \leq \frac{T^{-1/2} \mathbb{E}[\|\mathbf{c}\|_2^3] \mathbb{E}[\|\mathbf{A}\|_2^3] \mathbb{E}[\|\mathbf{u}_1\|_2^3]}{\{\mathbf{c}^\top \mathbb{E}[\mathbf{A} \text{Var}(\mathbf{u}_1) \mathbf{A}^\top] \mathbf{c}\}^{3/2}} \leq \frac{T^{-1/2} \mathbb{E}[\|\mathbf{A}\|_2^3] \mathbb{E}[\|\mathbf{u}_1\|_2^3]}{\{\max_i \text{Var}(u_{i1})\}^{3/2} \{\mathbb{E}[\|\mathbf{A}\|_{\mathbb{F}}^2]\}^{3/2}} \rightarrow 0$$

as T diverges to infinity. By the Lyapunov central limit theorem [Theorem 27.3 in [9], $\sum_{t=1}^T W_{nt}$ is hence asymptotically normal. Similarly, under Condition 6, we can show that \tilde{W}_{nt} is asymptotically normal. In addition, $\sum_{t=1}^T W_{nt}$ and $\sum_{t=1}^T \tilde{W}_{nt}$ are uncorrelated since $\{\mathbf{u}_t\}$ and $\{\mathbf{f}_t\}$ are uncorrelated mean zero processes. Therefore, $n^{1/2} T^{1/2} \mathbf{c}^\top (\tilde{\beta} - \beta)$ is asymptotically normal for any \mathbf{c} , and we have $\Sigma^{-1/2} (\tilde{\beta} - \beta) \xrightarrow{d} N(\mathbf{0}, \mathbf{I}_p)$, where $\Sigma = \mathbb{E}[(\mathbb{Z}_0^\top \mathbf{V}^{-1} \mathbb{Z}_0)^{-1}]$. By Theorem 2, we have $\sqrt{nT}(\tilde{\beta} - \beta) \xrightarrow{p} \mathbf{0}$. Notice that $\|\Sigma\|_{\mathbb{F}}^2 = O_p(nT)$, Slutsky's theorem yields $\Sigma^{-1/2} (\tilde{\beta} - \beta) \xrightarrow{d} N(\mathbf{0}, \mathbf{I}_p)$.

(ii) Similar to (i) and conditional on \mathbf{Z}_t and \mathbf{X} , the Lyapunov central limit theorem yields $(\mathbb{Z}_0^\top \mathbf{V}^{-1} \mathbb{Z}_0)^{1/2} (\tilde{\beta} - \beta) \xrightarrow{d} N(\mathbf{0}, \mathbf{I}_p)$, and Slutsky's theorem leads to $(\mathbb{Z}_0^\top \mathbf{V}^{-1} \mathbb{Z}_0)^{1/2} (\tilde{\beta} - \beta) \xrightarrow{d} N(\mathbf{0}, \mathbf{I}_p)$. \square

8.2. Technical results

We first collect some preliminary results for spline estimators in Lemmas 2 and 3.

Lemma 2. Under Condition 4, there exist constants c_1, c_2 such that $c_1 \|\mathbf{a}\|_2^2 \leq \mathbb{E}(\langle \mathbf{a}, \mathbf{a} \rangle_n) \leq c_2 \|\mathbf{a}\|_2^2$.

Proof. It follows from Condition 4 that, for any $\ell = 1, \dots, d$, the marginal density of X_ℓ on its support is bounded away from 0 and ∞ . Without loss of generality, we assume that, the support of \mathbf{X} is $[0, 1]^d$ and density $h(\mathbf{X})$ is bounded from below and above by m_1 and m_2 with $0 < m_1 \leq m_2 < \infty$.

Denote $f_\ell(X_\ell) = \sum_j a_{j\ell} \phi_j(X_\ell)$, $\ell \in \{1, \dots, d\}$ and $f_0 \equiv a_0$. Then, we have $\mathbb{E}(\langle \mathbf{a}, \mathbf{a} \rangle_n) = \mathbb{E}\{a_0 + \sum_j \sum_l a_{j\ell} \phi_j(X_\ell)\}^2 = \mathbb{E}\{a_0 + \sum_{\ell=1}^d f_\ell(X_\ell)\}^2$. Since the basis functions are centralized,

$$\mathbb{E}(\langle \mathbf{a}, \mathbf{a} \rangle_n) = \int_{\mathcal{X}} \{a_0 + \sum_{\ell=1}^d f_\ell(X_\ell)\}^2 h(\mathbf{X}) d\mathbf{X} \asymp \int_{\mathcal{X}} \{a_0 + \sum_{\ell=1}^d f_\ell(X_\ell)\}^2 d\mathbf{X} = a_0^2 + \int_{\mathcal{X}} \{\sum_{\ell=1}^d f_\ell(X_\ell)\}^2 d\mathbf{X}$$

By Lemma 1 of Stone [39], we obtain

$$\int_{\mathcal{X}} \{\sum_{\ell=1}^d f_\ell(X_\ell)\}^2 d\mathbf{X} \geq \left(\frac{C_0}{2}\right)^{d-1} \sum_{\ell=1}^d \int_0^1 f_\ell^2(x) dx = \left(\frac{C_0}{2}\right)^{d-1} (\sum_{\ell} \sum_j a_{j\ell}^2),$$

where $C_0 = 1 - (1 - m_1/m_2)^{1/2}$. Consequently, we have $\mathbb{E}(\langle \mathbf{a}, \mathbf{a} \rangle_n) \geq a_0^2 + (C_0/2)^{d-1} (\sum_{\ell} \sum_j a_{j\ell}^2) \geq \min\{1, (C_0/2)^d\} \|\mathbf{a}\|^2$. Similarly, we can establish that $\int_0^1 \{\sum_{\ell=1}^d f_\ell(X_\ell)\}^2 d\mathbf{X} \leq d^2 \sum_{\ell=1}^d \int_0^1 f_\ell^2(x) dx = d^2 (\sum_{\ell} \sum_j a_{j\ell}^2)$, and consequently, $\mathbb{E}(\langle \mathbf{a}, \mathbf{a} \rangle_n) \leq a_0^2 + d^2 (\sum_{\ell} \sum_j a_{j\ell}^2) \leq (1 + d^2) \|\mathbf{a}\|^2$. \square

Lemma 3. Under Condition 4, for some constant $C_1, C_2 > 0$, we have

$$\Pr \left\{ \sup_{\mathbf{a}_1, \mathbf{a}_2 \in \mathbb{R}^{jd+1}} \frac{|\langle \mathbf{a}_1, \mathbf{a}_2 \rangle_n - \mathbb{E}(\langle \mathbf{a}_1, \mathbf{a}_2 \rangle_n)|}{\sqrt{\mathbb{E}(\langle \mathbf{a}_1, \mathbf{a}_1 \rangle_n) \mathbb{E}(\langle \mathbf{a}_2, \mathbf{a}_2 \rangle_n)}} > s \right\} \leq C_1 J^2 \exp \left\{ -C_2 \frac{n}{J} \frac{s^2}{1+s} \right\}.$$

Proof. The proof is similar to that of Lemma A.2 in Huang et al. [23]. First, notice that $\langle \mathbf{a}_1, \mathbf{a}_2 \rangle_n - \mathbb{E}(\langle \mathbf{a}_1, \mathbf{a}_2 \rangle_n) = n^{-1} \sum_{\ell, \ell'} \sum_{j, j'} a_{j\ell}^{(1)} a_{j'\ell'}^{(2)} (\mathbf{e}_{j\ell}^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \mathbf{e}_{j'\ell'} - \mathbb{E}(\mathbf{e}_{j\ell}^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \mathbf{e}_{j'\ell'}))$, where $\mathbf{e}_{j\ell}$ is the $(J\ell + j + 1)$ th natural basis of \mathbb{R}^{jd+1} . Hence, we have $\boldsymbol{\Phi} \mathbf{e}_{j\ell} = \{\phi_j(X_{1\ell}), \dots, \phi_j(X_{n\ell})\}^\top$. For any j, j', ℓ, ℓ' , $\text{Var}(n^{-1} \mathbf{e}_{j\ell}^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \mathbf{e}_{j'\ell'}) \leq n^{-2} \sum_i \mathbb{E}(\phi_j^2(X_{i\ell}) \phi_{j'}^2(X_{i\ell'})) \lesssim n^{-1}$. As $|\phi_j(X_\ell)| \leq M$ for each j, ℓ for some $M > 0$,

$$\Pr \left\{ \left| \frac{1}{n} \mathbf{e}_{j\ell}^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \mathbf{e}_{j'\ell'} - \mathbb{E} \left(\frac{1}{n} \mathbf{e}_{j\ell}^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \mathbf{e}_{j'\ell'} \right) \right| > c_2 s \right\} \leq \exp \left\{ -\frac{ns^2}{M_1 + M_2 s} \right\}, \quad s > 0$$

by Bernstein's inequality with constants $M_1, M_2 > 0$. By the union bound,

$$\Pr \left[\bigcup_{j, j', \ell, \ell'} \left\{ \left| \frac{1}{n} \mathbf{e}_{j\ell}^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \mathbf{e}_{j'\ell'} - \mathbb{E} \left(\frac{1}{n} \mathbf{e}_{j\ell}^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \mathbf{e}_{j'\ell'} \right) \right| > \frac{c_2 s}{Jd} \right\} \right] \leq C_1 J^2 \exp \left\{ -\frac{C_2 ns^2}{J^2 + sJ} \right\}, \quad s > 0$$

for constants $C_1, C_2 > 0$. Denote $\mathcal{B} = \bigcup_{j, j', \ell, \ell'} \{ |n^{-1} \mathbf{e}_{j\ell}^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \mathbf{e}_{j'\ell'} - \mathbb{E}(n^{-1} \mathbf{e}_{j\ell}^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \mathbf{e}_{j'\ell'})| > c_2 s (Jd)^{-1} \}$, so that $\Pr(\mathcal{B}) < C_1 J^2 \exp \{ -C_2 ns^2 (J^2 + sJ)^{-1} \}$. For each $s > 0$, on \mathcal{B}^c , $|\langle \mathbf{a}_1, \mathbf{a}_2 \rangle_n - \mathbb{E}(\langle \mathbf{a}_1, \mathbf{a}_2 \rangle_n)| \leq \sum_{j, \ell} \sum_{j', \ell'} |a_{j\ell}^{(1)}| |a_{j'\ell'}^{(2)}| c_2 s / (Jd) \lesssim s \sqrt{\mathbb{E}(\langle \mathbf{a}_1, \mathbf{a}_1 \rangle_n) \mathbb{E}(\langle \mathbf{a}_2, \mathbf{a}_2 \rangle_n)}$, where the last inequality is due to Lemma 2. The conclusion follows. \square

Next, we document technical results for the proof of Theorem 1 in Lemmas 4–10 and Propositions 1–4.

Lemma 4. Under Conditions 3 and 5, for each n , with probability at least $1 - \delta$,

$$1 - n^{-1} \ln(1/\delta) \lesssim \lambda_{\min} \left(\frac{1}{n} \mathbf{G}^\top \mathbf{P} \mathbf{G} \right) < \lambda_{\max} \left(\frac{1}{n} \mathbf{G}^\top \mathbf{P} \mathbf{G} \right) \lesssim 1 + n^{-1} \ln(1/\delta).$$

Proof. Denote $\mathbf{R} = \mathbf{G} - \mathbf{P} \mathbf{G}$, and we have $\mathbf{G}^\top \mathbf{P} \mathbf{G} = \mathbf{G}^\top \mathbf{G} - \mathbf{G}^\top \mathbf{R}$. Thus, $\lambda_{\min} (n^{-1} \mathbf{G}^\top \mathbf{P} \mathbf{G}) \geq \lambda_{\min} (n^{-1} \mathbf{G}^\top \mathbf{G}) + \lambda_{\min} (-n^{-1} \mathbf{G}^\top \mathbf{R})$, and $\lambda_{\max} (n^{-1} \mathbf{G}^\top \mathbf{P} \mathbf{G}) \leq \lambda_{\max} (n^{-1} \mathbf{G}^\top \mathbf{G}) + \lambda_{\max} (-n^{-1} \mathbf{G}^\top \mathbf{R})$. Note that $\|\mathbf{R}\|_{\mathbb{F}}^2 \lesssim nJ^{-\kappa}$ by Condition 5. Thus, combining Condition 3, it holds that, with probability at least $1 - \delta$,

$$\|n^{-1} \mathbf{G}^\top \mathbf{R}\|_{\mathbb{F}}^2 = \frac{1}{n^2} \text{tr}(\mathbf{R}^\top \mathbf{G} \mathbf{G}^\top \mathbf{R}) \leq \lambda_{\max} \left(\frac{1}{n} \mathbf{G} \mathbf{G}^\top \right) \frac{1}{n} \text{tr}(\mathbf{R}^\top \mathbf{R}) \lesssim J^{-\kappa} \{1 + n^{-1} \ln(1/\delta)\}$$

and $|\lambda(\mathbf{G}^\top \mathbf{R}/n)| \lesssim J^{-\kappa} \{1 + n^{-1} \ln(1/\delta)\}$. By Condition 3, with probability at least $1 - \delta$, $1 - n^{-1} \ln(1/\delta) \lesssim \lambda_{\min} (n^{-1} \mathbf{G}^\top \mathbf{G}) < \lambda_{\max} (n^{-1} \mathbf{G}^\top \mathbf{G}) \lesssim 1 + n^{-1} \ln(1/\delta)$. The conclusion follows. \square

Lemma 5. Consider $\widehat{\boldsymbol{\beta}}^0$ satisfying $\|\widehat{\boldsymbol{\beta}}^0 - \boldsymbol{\beta}\|_2 = O_p(n^{-1/2+\alpha} T^{-1/2})$ for $\alpha \in [0, 1/2)$, such as the estimator in (14) in Section 8.3. Under Conditions 1 and 4–7, for $\widetilde{\mathbf{U}} = \mathbf{U} + \mathbb{Z}\{\mathbf{I}_T \otimes (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}^0)\}$ defined in Section 2.2,

- (i) $\mathbb{E}(\|\mathbf{F}^\top \widetilde{\mathbf{U}}\|_{\mathbb{F}}^2) = O((n + p^2 n^{2\alpha})T)$, $\mathbb{E}(\|\widetilde{\mathbf{U}}^\top \boldsymbol{\Phi}\|_{\mathbb{F}}^2) = O(nJ(T + p^2 n^{2\alpha}))$, $\mathbb{E}(\|\boldsymbol{\Phi}^\top \widetilde{\mathbf{U}} \mathbf{F}\|_{\mathbb{F}}^2) = O(p^2 n^{1+2\alpha} TJ)$, $\mathbb{E}(\|\widetilde{\mathbf{U}}^\top \boldsymbol{\Phi} \mathbf{B}\|_{\mathbb{F}}^2) \lesssim n(T + p^2 n^{2\alpha})$, and $\mathbb{E}(\|\mathbf{B}^\top \boldsymbol{\Phi}^\top \widetilde{\mathbf{U}} \mathbf{F}\|_{\mathbb{F}}^2) \lesssim p^2 n^{1+2\alpha} T$.
- (ii) With probability at least $1 - 5\delta$, $\|\mathbf{F}^\top \widetilde{\mathbf{U}}\|_{\mathbb{F}} \lesssim \{(n + p^2)T\}^{1/2} \{1 + \sqrt{\ln(1/\delta)}\}$, $\|\widetilde{\mathbf{U}}^\top \boldsymbol{\Phi}\|_{\mathbb{F}} \lesssim \{nJ(T + p^2)\}^{1/2} \{1 + \sqrt{\ln(1/\delta)}\}$, $\|\boldsymbol{\Phi}^\top \widetilde{\mathbf{U}} \mathbf{F}\|_{\mathbb{F}} \lesssim (p^2 n TJ)^{1/2} \{1 + \sqrt{\ln(1/\delta)}\}$, $\|\widetilde{\mathbf{U}}^\top \boldsymbol{\Phi} \mathbf{B}\|_{\mathbb{F}} \lesssim \sqrt{n(T + p^2 n^{2\alpha})} \{1 + \sqrt{\ln(1/\delta)}\}$, and $\|\mathbf{B}^\top \boldsymbol{\Phi}^\top \widetilde{\mathbf{U}} \mathbf{F}\|_{\mathbb{F}} \lesssim p \sqrt{n^{1+2\alpha} T} \{1 + \sqrt{\ln(1/\delta)}\}^2 \{1 + \sqrt{\ln(1/\delta)}\}$.
- (iii) With probability at least $1 - 4\delta$, $\|\widetilde{\mathbf{P}}\|_{\mathbb{F}} \lesssim \sqrt{J(T + p^2 n^{2\alpha})} \{1 + n^{-1} J \ln(J/\delta)\}^{3/2} \{1 + \sqrt{\ln(1/\delta)}\}$.

Proof.

- (i) By Lemma B.1 of [19], $\mathbb{E}(\|\mathbf{F}^\top \mathbf{U}\|_{\mathbb{F}}^2) = O(nT)$, $\mathbb{E}(\|\mathbf{U}^\top \boldsymbol{\Phi}\|_{\mathbb{F}}^2) = O(nJT)$, $\mathbb{E}(\|\boldsymbol{\Phi}^\top \mathbf{U} \mathbf{F}\|_{\mathbb{F}}^2) = O(nTJ)$, and $\mathbb{E}(\|\mathbf{P} \mathbf{U}\|_{\mathbb{F}}^2) = O(JT)$, and by Lemma C.6 in Fan et al. [19], $\mathbb{E}(\|\mathbf{U} \boldsymbol{\Phi} \mathbf{B}\|_{\mathbb{F}}^2) = O(nT)$ and $\mathbb{E}(\|\mathbf{B} \boldsymbol{\Phi}^\top \mathbf{U} \mathbf{F}\|_{\mathbb{F}}^2) = O(nT)$. Thus, it suffices to show

$$\mathbb{E}[\|\mathbb{Z}\{\mathbf{I}_T \otimes (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}^0)\} \mathbf{F}\|_{\mathbb{F}}^2] = O(p^2 T), \quad (9)$$

$$\mathbb{E}[\|\boldsymbol{\Phi}^\top \mathbb{Z}\{\mathbf{I}_T \otimes (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}^0)\}\|_{\mathbb{F}}^2] = O(p^2 nJ), \quad (10)$$

$$\mathbb{E}[\|\boldsymbol{\Phi}^\top \mathbb{Z}\{\mathbf{I}_T \otimes (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}^0)\} \mathbf{F}\|_{\mathbb{F}}^2] = O(p^2 nTJ), \quad (11)$$

$$\mathbb{E}[\|\mathbb{Z}\{\mathbf{I}_T \otimes (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}^0)\}^\top \boldsymbol{\Phi} \mathbf{B}\|_{\mathbb{F}}^2] = O(p^2 n^{1+2\alpha}), \quad (12)$$

$$E[\|\mathbf{B}^\top \Phi^\top \mathbb{Z}(\mathbf{I}_T \otimes (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^0))\mathbf{F}\|_{\mathbb{F}}^2] = O(p^2 n^{1+2\alpha} T). \quad (13)$$

By Proposition 5, $E[\|\mathbb{Z}(\mathbf{I}_T \otimes (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^0))\|_{\mathbb{F}}^2] \leq E[\|\mathbb{Z}\|_{\mathbb{F}}^2] \|\mathbf{I}_T\|_2^2 E[\|\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^0\|_{\mathbb{F}}^2] = O(p^2 n^{2\alpha})$. Then (9) follows from Cauchy-Schwarz inequality that $E[\|\mathbb{Z}(\mathbf{I}_T \otimes (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^0))\mathbf{F}\|_{\mathbb{F}}^2] \leq E[\|\mathbb{Z}(\mathbf{I}_T \otimes (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^0))\|_{\mathbb{F}}^2] E[\|\mathbf{F}\|_{\mathbb{F}}^2] = O(p^2 n^{2\alpha} T)$. As a consequence of Lemma 2, we have $E(\|\Phi\|_2^2) = O(n)$, and consequently $E(\|\Phi\|_{\mathbb{F}}^2) \leq (Jd+1)E(\|\Phi\|_2^2) = O(nJ)$, and (10) holds since $E[\|\Phi^\top \mathbb{Z}(\mathbf{I}_T \otimes (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^0))\|_{\mathbb{F}}^2] \leq E[\|\mathbb{Z}(\mathbf{I}_T \otimes (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^0))\|_{\mathbb{F}}^2] E[\|\Phi\|_{\mathbb{F}}^2] = O(p^2 n^{1+2\alpha} J)$. Applying Cauchy-Schwarz inequality, (11) follows $E[\|\Phi^\top \mathbb{Z}(\mathbf{I}_T \otimes (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^0))\mathbf{F}\|_{\mathbb{F}}^2] \leq E[\|\mathbb{Z}(\mathbf{I}_T \otimes (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^0))\|_{\mathbb{F}}^2] E[\|\Phi\|_{\mathbb{F}}^2] E[\|\mathbf{F}\|_{\mathbb{F}}^2] = O(p^2 n^{1+2\alpha} TJ)$. Also, (12) and (13) follow $E[\|\mathbb{Z}(\mathbf{I}_T \otimes (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^0))^\top \Phi \mathbf{B}\|_{\mathbb{F}}^2] \leq E[\|\mathbb{Z}(\mathbf{I}_T \otimes (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^0))\|_{\mathbb{F}}^2] E[\|\Phi\|_{\mathbb{F}}^2] + E[\|\mathbb{Z}(\mathbf{I}_T \otimes (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^0))\mathbf{R}\|_{\mathbb{F}}^2] = O(p^2 n^{1+2\alpha})$, and $E[\|\mathbf{B}^\top \Phi^\top \mathbb{Z}(\mathbf{I}_T \otimes (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^0))\mathbf{F}\|_{\mathbb{F}}^2] \leq E[\|\mathbf{G}^\top \mathbb{Z}(\mathbf{I}_T \otimes (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^0))\mathbf{F}\|_{\mathbb{F}}^2] + E[\|\mathbf{R}^\top \mathbb{Z}(\mathbf{I}_T \otimes (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^0))\mathbf{F}\|_{\mathbb{F}}^2] = O(p^2 n^{1+2\alpha} T)$, respectively.

(ii) By Condition 7 (iii), for $r_6 \in (0, 1)$ and any $x \geq 1$,

$$\tau(x) \leq 4 \max(b_1, b_2)^2 r_6 \{\max(r_2, r_3)(1 - r_6)/2\}^{2/\min(r_2, r_3)} \exp[2/\{\max(r_2, r_3)(1 - r_6)\}]\{2\alpha(x)\}^{r_6},$$

which implies that (f_{ik}, u_{it}) is τ -mixing [32] by Condition 7 (ii). Then, following Theorem 1 in Merlevède et al. [32] and Davydov's inequality [Corollary 16.2.4 in [1]], for each $k \in \{1, \dots, K\}$ and $i \in \{1, \dots, n\}$ and any $x > 0$, $\Pr(|\sum_{t=1}^T f_{ik} \tilde{u}_{it}| > x) \leq T \exp(-Cx^{r_{fu}}) + \exp\{-Cx^2/E(\sum_{t,s=1}^T f_{ik} \tilde{u}_{it} f_{sk} \tilde{u}_{is})\}$, where $r_{fu} = \{r_1^{-1} + \min(r_2, r_3)^{-1}\}^{-1}$. Let e^{-s} be the maximum of $T \exp(-Cx^{r_{fu}})$ and $\exp\{-Cx^2/E(\sum_{t,s=1}^T f_{ik} \tilde{u}_{it} f_{sk} \tilde{u}_{is})\}$. By Bonferroni inequality, with probability at least $1 - 2e^{-s}$, $|\sum_{t=1}^T f_{ik} \tilde{u}_{it}| \lesssim \max\{[s + \ln(pT)]^{1/r_{fu}}, \{E(\sum_{t,s=1}^T f_{ik} \tilde{u}_{it} f_{sk} \tilde{u}_{is})\}^{1/2} \sqrt{s + \ln p}\}$ uniformly for each $k \in \{1, \dots, K\}$ and $i \in \{1, \dots, n\}$. Then, by (i), with probability at least $1 - \delta$, $\|\mathbf{F}^\top \tilde{\mathbf{U}}\|_{\mathbb{F}} \lesssim \sqrt{(n + p^2 n^{2\alpha})T(s + \ln p)}$. The remaining two bounds follows similarly.

(iii) By Lemma 1, with probability at least $1 - \delta$, $\|\Phi\|_2^2 = \lambda_{\max}(\Phi^\top \Phi) \lesssim n\{1 + n^{-1}J \ln(J^2/\delta)\}$ and $\|(\Phi^\top \Phi)^{-1}\|_2 = \lambda_{\min}^{-1}(\Phi^\top \Phi) \lesssim n^{-1}\{1 + n^{-1}J \ln(J^2/\delta)\}$. Hence, with probability at least $1 - 4\delta$,

$$\|\tilde{\mathbf{P}}\tilde{\mathbf{U}}\|_{\mathbb{F}} \leq \|\Phi\|_2 \|(\Phi^\top \Phi)^{-1}\|_2 \|\Phi^\top \tilde{\mathbf{U}}\|_{\mathbb{F}} \lesssim \sqrt{(T + p^2 n^{2\alpha})\{1 + Jn^{-1} \ln(J^2/\delta)\}}^{3/2} \{1 + \sqrt{\ln(1/\delta)}\}. \quad \square$$

Denote \mathbf{K} a $K \times K$ diagonal matrix whose diagonals are the first K eigenvalues of $(nT)^{-1} \tilde{\mathbf{Y}}^\top \tilde{\mathbf{P}} \tilde{\mathbf{Y}}$. Then $(nT)^{-1} \tilde{\mathbf{Y}}^\top \tilde{\mathbf{P}} \tilde{\mathbf{Y}} \mathbf{F} = \hat{\mathbf{F}} \mathbf{K}$. Let $\mathbf{H} = (nT)^{-1} \mathbf{B}^\top \Phi^\top \Phi \mathbf{B} \mathbf{F}^\top \hat{\mathbf{F}} \mathbf{K}^{-1}$. Using $\tilde{\mathbf{Y}} = (\Phi \mathbf{B} + \mathbf{R})\mathbf{F}^\top + \tilde{\mathbf{U}}$ from Section 2.2.1, we have $\hat{\mathbf{F}} - \mathbf{F} \mathbf{H} = (\sum_{i=1}^8 \mathbf{A}_i) \mathbf{K}^{-1}$ where $\mathbf{A}_1 = (nT)^{-1} \mathbf{F} \mathbf{B}^\top \Phi^\top \tilde{\mathbf{U}} \hat{\mathbf{F}}$, $\mathbf{A}_2 = (nT)^{-1} \tilde{\mathbf{U}}^\top \Phi \mathbf{B} \mathbf{F}^\top \hat{\mathbf{F}}$, $\mathbf{A}_3 = (nT)^{-1} \tilde{\mathbf{U}}^\top \tilde{\mathbf{P}} \tilde{\mathbf{U}} \hat{\mathbf{F}}$, $\mathbf{A}_4 = (nT)^{-1} \mathbf{F} \mathbf{B}^\top \Phi^\top \mathbf{R} \mathbf{F}^\top \hat{\mathbf{F}}$, $\mathbf{A}_5 = (nT)^{-1} \mathbf{F} \mathbf{R}^\top \Phi \mathbf{B} \mathbf{F}^\top \hat{\mathbf{F}}$, $\mathbf{A}_6 = (nT)^{-1} \mathbf{F} \mathbf{R}^\top \mathbf{P} \mathbf{R} \mathbf{F}^\top \hat{\mathbf{F}}$, $\mathbf{A}_7 = (nT)^{-1} \mathbf{F} \mathbf{R}^\top \tilde{\mathbf{P}} \tilde{\mathbf{U}} \hat{\mathbf{F}}$, and $\mathbf{A}_8 = (nT)^{-1} \tilde{\mathbf{U}}^\top \mathbf{P} \mathbf{R} \mathbf{F}^\top \hat{\mathbf{F}}$. Next, in Lemmas 6–10, we will provide a bound on $\|\mathbf{H} - \mathbf{I}\|_{\mathbb{F}}$ in probability.

Lemma 6. With probability at least $1 - 5\delta$, $\|\mathbf{K}^{-1}\|_2 \lesssim 1 + n^{-1} \ln(1/\delta)$.

Proof. The K largest eigenvalues of $(nT)^{-1} \tilde{\mathbf{Y}}^\top \tilde{\mathbf{P}} \tilde{\mathbf{Y}}$ are the same as those of $\mathbf{W} = (nT)^{-1} (\Phi^\top \Phi)^{-1/2} \Phi^\top \tilde{\mathbf{Y}} \tilde{\mathbf{Y}}^\top \Phi (\Phi^\top \Phi)^{-1/2}$. Substituting $\tilde{\mathbf{Y}} = \mathbf{G} \mathbf{F}^\top + \tilde{\mathbf{U}}$ and $T^{-1} \mathbf{F}^\top \mathbf{F} = \mathbf{I}_K$, we have $\mathbf{W} = \sum_{i=1}^4 \mathbf{W}_i$ where $\mathbf{W}_1 = n^{-1} (\Phi^\top \Phi)^{-1/2} \Phi^\top \mathbf{G} \mathbf{G}^\top \Phi (\Phi^\top \Phi)^{-1/2}$, $\mathbf{W}_2 = (nT)^{-1} (\Phi^\top \Phi)^{-1/2} \Phi^\top \mathbf{G} \mathbf{F}^\top \tilde{\mathbf{U}}^\top \Phi (\Phi^\top \Phi)^{-1/2}$, $\mathbf{W}_3 = \mathbf{W}_2^\top$, and $\mathbf{W}_4 = (nT)^{-1} (\Phi^\top \Phi)^{-1/2} \Phi^\top \tilde{\mathbf{U}} \tilde{\mathbf{U}}^\top \Phi (\Phi^\top \Phi)^{-1/2}$. By Lemma 1, with probability at least $1 - \delta$, $\|\Phi\|_2^2 = \lambda_{\max}(\Phi^\top \Phi) \lesssim n\{1 + n^{-1}J \ln(J^2/\delta)\}$ and $\|(\Phi^\top \Phi)^{-1}\|_2 = \lambda_{\min}^{-1}(\Phi^\top \Phi) \lesssim n^{-1}\{1 + n^{-1}J \ln(J^2/\delta)\}$. By Lemma 4, with probability at least $1 - \delta$, $\|\mathbf{P} \mathbf{G}\|_2^2 = \lambda_{\max}(\mathbf{G}^\top \mathbf{P} \mathbf{G}) \lesssim n(1 + J^{-K})\{1 + n^{-1} \ln(1/\delta)\}$. Hence, with probability at least $1 - 5\delta$,

$$\begin{aligned} \|\mathbf{W}_2\|_2 &\leq \frac{1}{n} \|(\Phi^\top \Phi)^{-1/2}\|_2^2 \|\Phi\|_2 \|\mathbf{P} \mathbf{G}\|_2 \left\| \frac{1}{T} \mathbf{F}^\top \tilde{\mathbf{U}}^\top \Phi \right\|_{\mathbb{F}} \\ &\lesssim p \sqrt{J n^{2\alpha-1} T^{-1} (1 + J^{-K}) \{1 + n^{-1}J \ln(J^2/\delta)\}^{3/2} \{1 + \sqrt{\ln(1/\delta)}\} \{1 + n^{-1} \ln(1/\delta)\}}, \end{aligned}$$

and by Lemma 5, with probability at least $1 - 4\delta$,

$$\|\mathbf{W}_4\|_2 \leq \frac{1}{nT} \|(\Phi^\top \Phi)^{-1/2}\|_2^2 \|\Phi^\top \tilde{\mathbf{U}}\|_{\mathbb{F}}^2 \lesssim \frac{J(T + p^2 n^{2\alpha})}{nT} \{1 + J \ln(J^2/\delta)/n\} \{1 + \sqrt{\ln(1/\delta)}\}.$$

By Weyl's Theorem, $|\lambda_k(\mathbf{W}) - \lambda_k(\mathbf{W}_1)| \leq \|\mathbf{W} - \mathbf{W}_1\|_2$ for each $k \in \{1, \dots, K\}$. Hence, with probability at least $1 - 5\delta$,

$$|\lambda_k(\mathbf{W}) - \lambda_k(\mathbf{W}_1)| \lesssim [p \sqrt{J(n^{1-2\alpha} T)^{-1/2} + \{J(T + p^2 n^{2\alpha})\}/(nT)}] \{1 + J \ln(J^2/\delta)/n\}^{3/2} \{1 + \sqrt{\ln(1/\delta)}\}.$$

Note that the K largest eigenvalues of \mathbf{W}_1 is also the K largest eigenvalues of $n^{-1} \mathbf{G}^\top \mathbf{P} \mathbf{G}$. Thus, by Lemma 4, with probability at least $1 - 5\delta$, $\|\mathbf{K}^{-1}\|_2 \lesssim 1 + n^{-1} \ln(1/\delta)$. \square

Lemma 7. With probability at least $1 - 7\delta$, (i) $\|\mathbf{A}_1\|_{\mathbb{F}}, \|\mathbf{A}_2\|_{\mathbb{F}} \lesssim \sqrt{n^{-1}(T + p^2 n^{2\alpha})\{1 + \sqrt{\ln(1/\delta)}\}}$, (ii) $\|\mathbf{A}_3\|_{\mathbb{F}} \lesssim n^{-1} T^{-1/2} J(T + p^2 n^{2\alpha})\{1 + \sqrt{\ln(1/\delta)}\}$, (iii) $\|\mathbf{A}_4\|_{\mathbb{F}}, \|\mathbf{A}_5\|_{\mathbb{F}} \lesssim (J^{-K/2} \sqrt{T})\{1 + \sqrt{\ln(1/\delta)}\}$, (iv) $\|\mathbf{A}_7\|_{\mathbb{F}}, \|\mathbf{A}_8\|_{\mathbb{F}} \lesssim \sqrt{(T + p^2 n^{2\alpha})(nJ^{K-1})^{-1}} \{1 + \sqrt{\ln(1/\delta)}\}$; and $\|\mathbf{A}_6\|_{\mathbb{F}} \lesssim J^{-K} \sqrt{T}$.

Proof. Notice that $\|\mathbf{F}\|_{\mathbb{F}} = \sqrt{KT}$ with probability 1 and $\|\widehat{\mathbf{F}}\|_{\mathbb{F}} = \sqrt{KT}$. Then, both (i) and (ii) follow from Lemma 5. By Condition 5 that $\|\mathbf{R}\|_{\mathbb{F}}^2 \lesssim nJ^{-\kappa}$, (iii) follows from Lemma 4 and $\Phi\mathbf{B} = \mathbf{P}\mathbf{G}$. Part (iv) follows from Lemma 5 and $\|\mathbf{R}\|_{\mathbb{F}}^2 \lesssim nJ^{-\kappa}$. Result on \mathbf{A}_6 follows similarly to (iii) given $\|\mathbf{P}\|_2 = 1$. \square

Lemma 8. With probability at least $1 - 3\delta$, (i) $\|\mathbf{A}_1\|_{\max}, \|\mathbf{A}_2\|_{\max} \lesssim n^{-1/2}T^{-1}\sqrt{T + p^2n^{2\alpha}\{\ln(T)\}^{2/r_2}\{1 + \ln(1/\delta)\}}$, (ii) $\|\mathbf{A}_3\|_{\max} \lesssim n^{-1/2}T^{-1}\sqrt{T + p^2n^{2\alpha}\{\ln(T)\}^{1/r_2}\{1 + \ln(1/\delta)\}}$, (iii) $\|\mathbf{A}_4\|_{\max}, \|\mathbf{A}_5\|_{\max} \lesssim n^{-1}T^{-1}\{\ln(T)\}^{3/r_2}J^{-\kappa}\{1 + \ln(1/\delta)\}$, (iv) $\|\mathbf{A}_7\|_{\max}, \|\mathbf{A}_8\|_{\max} \lesssim (nT)^{-1}J^{-\kappa}\sqrt{J(T + p^2n^{2\alpha}\{\ln(T)\}^{2/r_2}\{1 + \ln(1/\delta)\})}$; and $\|\mathbf{A}_6\|_{\max} \lesssim (nT)^{-1}J^{-2\kappa}\{\ln(T)\}^{3/r_2}$.

Proof. By Lemma B.1 in Fan et al. [17], with probability at least $1 - \delta$, $\|\widetilde{\mathbf{U}}\widetilde{\mathbf{P}}\widetilde{\mathbf{U}}\|_{\max} \lesssim \sqrt{n}(T + p^2n^{2\alpha})\{1 + \ln(1/\delta)\}$. Also, the proof of Lemma D.2 in Wang and Fan [43] implies that $\|\mathbf{U}^{\top}\Phi\mathbf{B}\|_{\infty} \lesssim \sqrt{n}T$. Hence, with probability at least $1 - \delta$, $\|\widetilde{\mathbf{U}}^{\top}\Phi\mathbf{B}\|_{\infty} \lesssim \sqrt{n}(T + p^2n^{2\alpha})\{1 + \ln(1/\delta)\}$ by Lemma 5. Then, the results follow from that $\|\mathbf{F}\|_{\max} \lesssim \{\ln(T) + \ln(1/\delta)\}^{1/r_2}$ with probability at least $1 - \delta$. \square

Proposition 1. Given $Jd + 1 < n$ and $\kappa \geq 1$,

- (i) With probability at least $1 - 12\delta$, $T^{-1}\|\widehat{\mathbf{F}} - \mathbf{F}\mathbf{H}\|_{\mathbb{F}}^2 \lesssim (n^{-1} + \{n^{1-2\alpha}T\}^{-1}p^2 + J^{-\kappa})\{1 + \sqrt{\ln(1/\delta)}\}^2\{1 + n^{-1}\ln(1/\delta)\}$.
- (ii) With probability at least $1 - 8\delta$, $\|\widehat{\mathbf{F}} - \mathbf{F}\mathbf{H}\|_{\max} \lesssim (n^{-1/2} + \{\sqrt{n^{1-2\alpha}T}\}^{-1}p)\{\ln(T)\}^{2/r_2}\{1 + \ln(1/\delta)\}$.

Proof. By Lemma 6, $\|\mathbf{K}^{-1}\|_2 \lesssim 1 + n^{-1}\ln(1/\delta)$ with probability at least $1 - 5\delta$. The result follows from Lemmas 7 and 8. \square

Lemma 9. With probability at least $1 - 20\delta$, (i) $T^{-1}\|\mathbf{A}_1\|_{\mathbb{F}}^2 \lesssim \{n^{-2} + n^{-1+2\alpha}T^{-1}p^2 + (nTJ^{\kappa})^{-1}(T + p^2n^{2\alpha})\}\{1 + \sqrt{\ln(1/\delta)}\}^2\{1 + n^{-1}\ln(1/\delta)\}$, (ii) $T^{-2}\|\mathbf{F}^{\top}\mathbf{A}_2\|_{\mathbb{F}}^2 \lesssim n^{-1+2\alpha}T^{-1}p^2\{1 + \sqrt{\ln(1/\delta)}\}^2$, (iii) $T^{-2}\|\mathbf{F}^{\top}(\widehat{\mathbf{F}} - \mathbf{F}\mathbf{H})\|_{\mathbb{F}}^2 \lesssim \{n^{-2} + n^{-1+2\alpha}T^{-1}p^2 + J^{-\kappa}\}\{1 + \sqrt{\ln(1/\delta)}\}^2$, and (iv) $T^{-2}\|\widehat{\mathbf{F}}^{\top}(\widehat{\mathbf{F}} - \mathbf{F}\mathbf{H})\|_{\mathbb{F}}^2 \lesssim \{n^{-2} + n^{-1+2\alpha}T^{-1}p^2 + J^{-\kappa}\}\{1 + \sqrt{\ln(1/\delta)}\}^2$.

Proof. (i) First, by Lemmas 4 and 6, with probability at least $1 - 6\delta$, $\|\mathbf{H}\|_2 \leq (nT)^{-1}\|\mathbf{P}\mathbf{G}\|_{\mathbb{F}}^2\|\mathbf{F}\|_{\mathbb{F}}\|\widehat{\mathbf{F}}\|_{\mathbb{F}}\|\mathbf{K}^{-1}\|_2 \lesssim 1 + n^{-1}\ln(1/\delta)$. Then, by Lemma 5 and Proposition 1, with probability at least $1 - 20\delta$,

$$\begin{aligned} \|\mathbf{B}^{\top}\Phi^{\top}\widetilde{\mathbf{U}}\widetilde{\mathbf{F}}\|_{\mathbb{F}}^2 &\leq 2\|\mathbf{B}^{\top}\Phi^{\top}\widetilde{\mathbf{U}}(\widehat{\mathbf{F}} - \mathbf{F}\mathbf{H})\|_{\mathbb{F}}^2 + 2\|\mathbf{B}^{\top}\Phi^{\top}\widetilde{\mathbf{U}}\mathbf{F}\mathbf{H}\|_{\mathbb{F}}^2 \\ &\lesssim \{T^2 + p^2n^{2\alpha}T + p^4n^{4\alpha} + nT(T + p^2n^{2\alpha})/J^{\kappa}\}\{1 + \sqrt{\ln(1/\delta)}\}^2\{1 + n^{-1}\ln(1/\delta)\}. \end{aligned}$$

The result follows that $\|\mathbf{F}\|_{\mathbb{F}} = \|\widehat{\mathbf{F}}\|_{\mathbb{F}} = \sqrt{KT}$ with probability 1.

(ii) By Lemma 5, with probability at least $1 - 4\delta$, $T^{-2}\|\mathbf{F}^{\top}\mathbf{A}_2\|_{\mathbb{F}}^2 \leq n^{-2}T^{-4}\|\mathbf{F}^{\top}\widetilde{\mathbf{U}}^{\top}\Phi\mathbf{B}\|_{\mathbb{F}}^2\|\mathbf{F}\|_{\mathbb{F}}^2\|\widehat{\mathbf{F}}\|_{\mathbb{F}}^2 \lesssim (nT)^{-1}p^2\{1 + \sqrt{\ln(1/\delta)}\}^2$.

(iii) Combining (i) and (ii), the result follows from Lemma 7.

(iv) The result follows from $T^{-1}\|\widehat{\mathbf{F}}^{\top}(\widehat{\mathbf{F}} - \mathbf{F}\mathbf{H})\|_{\mathbb{F}} \leq T^{-1}\|\widehat{\mathbf{F}} - \mathbf{F}\mathbf{H}\|_{\mathbb{F}}^2 + T^{-1}\|\mathbf{H}^{\top}\mathbf{F}^{\top}(\widehat{\mathbf{F}} - \mathbf{F}\mathbf{H})\|_{\mathbb{F}}$. \square

Lemma 10. With probability at least $1 - 20\delta$,

$$\|\mathbf{H}^{\top}\mathbf{H} - \mathbf{I}_K\|_{\mathbb{F}}^2 \lesssim \left(\frac{1}{n^2} + \frac{p^2}{n^{1-2\alpha}T} + \frac{1}{J^{\kappa}}\right)\{1 + \sqrt{\ln(1/\delta)}\}^2\{1 + n^{-1}\ln(1/\delta)\}.$$

Proof. By Condition 2, $\mathbf{F}^{\top}\mathbf{F} = T\mathbf{I}_K$ with probability 1 and $\widehat{\mathbf{F}}^{\top}\widehat{\mathbf{F}} = T\mathbf{I}_K$. So $\mathbf{H}^{\top}\mathbf{H} = T^{-1}(\mathbf{F}\mathbf{H})^{\top}\mathbf{F}\mathbf{H} = T^{-1}(\mathbf{F}\mathbf{H} - \widehat{\mathbf{F}})^{\top}\mathbf{F}\mathbf{H} + T^{-1}\widehat{\mathbf{F}}^{\top}(\mathbf{F}\mathbf{H} - \widehat{\mathbf{F}}) + \mathbf{I}_K$ and $\|\mathbf{H}^{\top}\mathbf{H} - \mathbf{I}_K\|_{\mathbb{F}} \leq T^{-1}\|(\mathbf{F}\mathbf{H} - \widehat{\mathbf{F}})^{\top}\mathbf{F}\|_{\mathbb{F}}\|\mathbf{H}\|_2 + T^{-1}\|\mathbf{F}^{\top}(\widehat{\mathbf{F}} - \mathbf{F}\mathbf{H})\|_{\mathbb{F}}$, which gives the desired result. \square

Define $\widehat{\mathbf{B}} = T^{-1}(\Phi^{\top}\Phi)^{-1}\Phi^{\top}\widetilde{\mathbf{Y}}\widehat{\mathbf{F}}$ so that $\widehat{\mathbf{G}} = T^{-1}\mathbf{P}\widetilde{\mathbf{Y}}\widehat{\mathbf{F}} = \Phi\widehat{\mathbf{B}}$, we have $\widehat{\mathbf{B}} - \mathbf{B}\mathbf{H} = \sum_{i=1}^4\mathbf{C}_i$ where $\mathbf{C}_1 = T^{-1}(\Phi^{\top}\Phi)^{-1}\Phi^{\top}\mathbf{R}\mathbf{F}^{\top}\widehat{\mathbf{F}}$, $\mathbf{C}_2 = T^{-1}(\Phi^{\top}\Phi)^{-1}\Phi^{\top}\widetilde{\mathbf{U}}\mathbf{F}\mathbf{H}$, $\mathbf{C}_3 = T^{-1}(\Phi^{\top}\Phi)^{-1}\Phi^{\top}\widetilde{\mathbf{U}}(\widehat{\mathbf{F}} - \mathbf{F}\mathbf{H})$, and $\mathbf{C}_4 = T^{-1}\mathbf{B}\mathbf{F}^{\top}(\widehat{\mathbf{F}} - \mathbf{F}\mathbf{H})$.

Proposition 2. With probability at least $1 - 20\delta$,

- (i) $\|\widehat{\mathbf{B}} - \mathbf{B}\mathbf{H}\|_{\mathbb{F}}^2 \lesssim \{n^{-2}J + n^{-1+2\alpha}T^{-1}p^2J + n^{-2+4\alpha}T^{-2}p^4J + J^{-\kappa+1}\}\{1 + J\ln(J^2/\delta)/n\}^3\{1 + \sqrt{\ln(1/\delta)}\}^4$,
- (ii) $n^{-1}\|\widehat{\mathbf{G}} - \mathbf{G}\mathbf{H}\|_{\mathbb{F}}^2 \lesssim \{n^{-2}J + n^{-1+2\alpha}T^{-1}p^2J + n^{-2+4\alpha}T^{-2}p^4J + J^{-\kappa+1}\}\{1 + J\ln(J^2/\delta)/n\}^4\{1 + \sqrt{\ln(1/\delta)}\}^4$.

Proof. (i) By Lemmas 1, 5 and 9, with probability at least $1 - 20\delta$,

$$\|\mathbf{C}_1\|_{\mathbb{F}}^2 \lesssim \frac{1}{J^{\kappa}}\{1 + J\ln(J^2/\delta)/n\}^3,$$

$$\|\mathbf{C}_2\|_{\mathbb{F}}^2 \lesssim \frac{p^2J}{n^{2\alpha}T}\{1 + J\ln(J^2/\delta)/n\}^2\{1 + \sqrt{\ln(1/\delta)}\}^2,$$

$$\|\mathbf{C}_3\|_{\mathbb{F}}^2 \lesssim \left(\frac{J}{n^2} + \frac{p^2 J}{n^{2-2\alpha} T} + \frac{p^4 J}{n^{2-4\alpha} T^2} + \frac{T + p^2}{n T^{\kappa-1}} \right) \{1 + J \ln(J^2/\delta)/n\}^2 \{1 + \sqrt{\ln(1/\delta)}\}^4,$$

$$\|\mathbf{C}_4\|_{\mathbb{F}}^2 \lesssim \left(\frac{J}{n^2} + \frac{p^2 J}{n^{1-2\alpha} T} + \frac{1}{J^{\kappa-1}} \right) \{1 + J \ln(J^2/\delta)/n\}^3 \{1 + \sqrt{\ln(1/\delta)}\}^2.$$

So $\|\widehat{\mathbf{B}} - \mathbf{B}\mathbf{H}\|_{\mathbb{F}}^2 \lesssim \{n^{-2}J + n^{-1+2\alpha}T^{-1}p^2J + n^{-2+4\alpha}T^{-2}p^4J + J^{-\kappa+1}\} \{1 + J \ln(J^2/\delta)/n\}^3 \{1 + \sqrt{\ln(1/\delta)}\}^4$.

(ii) The result follows from $n^{-1}\|\widehat{\mathbf{G}} - \mathbf{G}\mathbf{H}\|_{\mathbb{F}}^2 \leq n^{-2}\|\Phi(\widehat{\mathbf{B}} - \mathbf{B}\mathbf{H})\|_{\mathbb{F}}^2 + n^{-2}\|\mathbf{R}\mathbf{H}\|_{\mathbb{F}}^2$. \square

Proposition 3. With probability at least $1 - 20\delta$, (i) $\|\widehat{\mathbf{B}} - \mathbf{B}\mathbf{H}\|_{\max} \lesssim n^{-1/2}T^{-1}\{(T + p^2n^{2\alpha})\ln(n)\}^{1/2}\{1 + \ln(1/\delta)\}$, (ii) $\|\widehat{\mathbf{G}} - \mathbf{G}\mathbf{H}\|_{\max} \lesssim T^{-1}\{(T + p^2n^{2\alpha})\ln(n)\}^{1/2}\{1 + \ln(1/\delta)\}$, and (iii) $\|\widehat{\mathbf{G}} - \mathbf{G}\mathbf{H}^{-1}\|_{\max} \lesssim T^{-1}\{(T + p^2n^{2\alpha})\ln(n)\}^{1/2}\{1 + \ln(1/\delta)\}$.

Proof. (i) By Lemma B.1 in Fan et al. [17], with probability at least $1 - \delta$, $\|\widehat{\mathbf{F}}\mathbf{U}\|_{\max} \lesssim \sqrt{(T + p^2)\ln(n)\{1 + \ln(1/\delta)\}}$. Then, by Lemmas 1, 5, and 9, with probability at least $1 - 20\delta$, $\|\mathbf{C}_1\|_{\max} \lesssim \{\sqrt{nT}^{\kappa}\}^{-1}\{\ln(T)\}^{2/2}\{1 + \ln(1/\delta)\}$, $\|\mathbf{C}_2\|_{\max} \lesssim \{\sqrt{nT}\}^{-1}\sqrt{(T + p^2n^{2\alpha})\ln(n)\{1 + \ln(1/\delta)\}}$, $\|\mathbf{C}_3\|_{\max} \lesssim \{nT^2\}^{-1}(T + p^2n^{2\alpha})\{\ln(T)\}^{2/2}\{1 + \ln(1/\delta)\}$, and $\|\mathbf{C}_4\|_{\max} \lesssim \{\sqrt{nT}^2J^{\kappa}\}^{-1}(T + p^2n^{2\alpha})\{\ln(T)\}^{2/2}\{1 + \ln(1/\delta)\}\delta$. So $\|\widehat{\mathbf{B}} - \mathbf{B}\mathbf{H}\|_{\max} \lesssim n^{-1/2}T^{-1}\{(T + p^2n^{2\alpha})\ln(n)\}^{1/2}\{1 + \ln(1/\delta)\}$.

(ii) The result follows from $\|\widehat{\mathbf{G}} - \mathbf{G}\mathbf{H}\|_{\max} \leq n^{-2}\|\Phi(\widehat{\mathbf{B}} - \mathbf{B}\mathbf{H})\|_{\max} + n^{-2}\|\mathbf{R}\mathbf{H}\|_{\max}\delta$.

(iii) The result follows from $\widehat{\mathbf{G}} - \mathbf{G}\mathbf{H}^{-1} = T^{-1}\mathbf{G}\mathbf{H}^{-1}(\mathbf{H}\mathbf{F}^{\top} - \widehat{\mathbf{F}}^{\top})\widehat{\mathbf{F}} + T^{-1}\mathbf{P}\widehat{\mathbf{U}}(\widehat{\mathbf{F}} - \mathbf{F}\mathbf{H}) + T^{-1}\mathbf{P}\widehat{\mathbf{U}}\mathbf{F}\mathbf{H}$. \square

Proposition 4. With probability at least $1 - 20\delta$,

$$\|\mathbf{H} - \mathbf{I}_K\|_{\mathbb{F}}^2 \lesssim \left(\frac{1}{n^2} + \frac{p^2}{n^{1-2\alpha}T} + \frac{1}{J^{\kappa}} \right) \{1 + \sqrt{\ln(1/\delta)}\}^2 \{1 + n^{-1}\ln(1/\delta)\}.$$

Proof. Note that $\mathbf{H}\mathbf{K} = n^{-1}\mathbf{B}^{\top}\Phi^{\top}\Phi\mathbf{B}(T^{-1}\mathbf{F}^{\top}\widehat{\mathbf{F}} - \mathbf{H}) + n^{-1}\mathbf{B}^{\top}\Phi^{\top}\Phi\mathbf{B}\mathbf{H}$. By Lemma 9, with probability at least $1 - 20\delta$,

$$\|n^{-1}\mathbf{B}^{\top}\Phi^{\top}\Phi\mathbf{B}(T^{-1}\mathbf{F}^{\top}\widehat{\mathbf{F}} - \mathbf{H})\|_{\mathbb{F}} \lesssim \{n^{-1} + p(n^{1-2\alpha}T)^{-1/2} + J^{-\kappa/2}\} \{1 + \sqrt{\ln(1/\delta)}\} \sqrt{1 + n^{-1}\ln(1/\delta)}.$$

In addition, by Conditions 3 and 5, $\|\mathbf{G}^{\top}\mathbf{G} - \mathbf{B}^{\top}\Phi^{\top}\Phi\mathbf{B}\|_{\mathbb{F}} \lesssim nJ^{-\kappa/2}$. Therefore, with probability at least $1 - 20\delta$,

$$\|n^{-1}\mathbf{G}^{\top}\mathbf{G}\mathbf{H} - \mathbf{H}\mathbf{K}\|_{\mathbb{F}} \lesssim \{n^{-1} + p(n^{1-2\alpha}T)^{-1/2} + J^{-\kappa/2}\} \{1 + \sqrt{\ln(1/\delta)}\} \sqrt{1 + n^{-1}\ln(1/\delta)}.$$

This implies that with probability at least $1 - 20\delta$, \mathbf{H} (up to an error term) is a matrix consisting of eigenvectors of $n^{-1}\mathbf{G}^{\top}\mathbf{G}$. By Condition 2, $\mathbf{G}^{\top}\mathbf{G}$ is a diagonal matrix with distinct eigenvalues with probability 1. Thus, each eigenvalue is associated with a unique unitary eigenvector up to a sign change and each eigenvector has a single non-zero entry. Thus, with probability at least $1 - 20\delta$,

$$\|\mathbf{H} - \mathbf{D}\|_{\mathbb{F}} \lesssim \{n^{-1} + p(n^{1-2\alpha}T)^{-1/2} + J^{-\kappa/2}\} \{1 + \sqrt{\ln(1/\delta)}\} \sqrt{1 + n^{-1}\ln(1/\delta)}$$

for some diagonal matrix \mathbf{D} . By Lemma 10, with probability at least $1 - 20\delta$, for each $k \in \{1, \dots, K\}$,

$$|\lambda_k(\mathbf{H}) - \eta| \lesssim \{n^{-1} + p(n^{1-2\alpha}T)^{-1/2} + J^{-\kappa/2}\} \{1 + \sqrt{\ln(1/\delta)}\} \sqrt{1 + n^{-1}\ln(1/\delta)},$$

where η is either 1 or -1 . Without loss of generality, we can assume that all entries of \mathbf{H} is positive (otherwise we can multiply the corresponding columns of $\widehat{\mathbf{F}}$ and $\widehat{\mathbf{G}}$ by -1). Hence, with probability at least $1 - 20\delta$,

$$\|\mathbf{H} - \mathbf{I}_K\|_{\mathbb{F}}^2 = \sum_{i \neq j} h_{ij}^2 + \sum_{i=1}^K (h_{ii} - 1)^2 \lesssim \{n^{-2} + p^2(n^{1-2\alpha}T)^{-1} + J^{-\kappa}\} \{1 + \sqrt{\ln(1/\delta)}\}^2 \{1 + n^{-1}\ln(1/\delta)\}. \quad \square$$

Finally, we present technical results for establishing Theorem 2 in Lemmas 11–14. Recall that $\mathcal{V}(\mathbf{f}_t) = T^{-2} \sum_{t=-T+1}^{T-1} (T - |t|) \widehat{\Sigma}_{\mathbf{f}}(t)$ as defined in Section 2.2.2, where $\widehat{\Sigma}_{\mathbf{f}}(s) = (T - s)^{-1} \sum_{t=1}^{T-s} (\mathbf{f}_t - \bar{\mathbf{f}})(\mathbf{f}_{t+s} - \bar{\mathbf{f}})^{\top}$ and $\widehat{\Sigma}_{\mathbf{f}}(-s) = (T - s)^{-1} \sum_{t=s}^T (\mathbf{f}_t - \bar{\mathbf{f}})(\mathbf{f}_{t-s} - \bar{\mathbf{f}})^{\top}$ for $s \geq 0$, respectively.

Lemma 11. Under Condition 2, with probability at least $1 - \delta$,

$$\|\mathcal{V}(\widehat{\mathbf{f}}_t) - \mathcal{V}(\mathbf{f}_t)\|_{\mathbb{F}} \lesssim \frac{1}{T} \left(\frac{1}{\sqrt{n}} + \frac{p}{\sqrt{n^{1-2\alpha}T}} + \frac{1}{J^{\kappa/2}} \right) \{1 + \sqrt{\ln(20/\delta)}\}.$$

Proof.

Note that $\mathcal{V}(\mathbf{f}_t) = T^{-2} \sum_{t,s=1}^T (\mathbf{f}_t - \bar{\mathbf{f}})(\mathbf{f}_s - \bar{\mathbf{f}})^{\top} = T^{-2} \mathbf{F}^{\top} \mathbf{P}_1 \mathbf{F}$, where \mathbf{P}_1 is the projection matrix onto $(1, \dots, 1)^{\top} \in \mathbb{R}^T$. Thus, by Theorem 1

$$\|\mathcal{V}(\widehat{\mathbf{f}}_t) - \mathcal{V}(\mathbf{f}_t)\|_{\mathbb{F}}^2 = \frac{1}{T^4} \|\widehat{\mathbf{F}} \mathbf{P}_1 \widehat{\mathbf{F}}^{\top} - \mathbf{F} \mathbf{P}_1 \mathbf{F}^{\top}\|_{\mathbb{F}}^2 \lesssim \frac{1}{T^2} \left(\frac{1}{n} + \frac{p^2}{n^{1-2\alpha}T} + \frac{1}{J^{\kappa}} \right) \{1 + \sqrt{\ln(20/\delta)}\}^2.$$

The conclusion follows. \square

Lemma 12. Under Conditions 1, 2, and 7, $\|\mathcal{V}(\mathbf{f}_t) - \text{Var}(T^{-1} \sum_{t=1}^T \mathbf{f}_t)\|_{\mathbb{F}} \lesssim T^{-2}$.

Proof. Recall that $\text{Var}(T^{-1} \sum_{t=1}^T \mathbf{f}_t) = T^{-2} \sum_{t,s} \text{Cov}(\mathbf{f}_t, \mathbf{f}_s)$ and $\mathcal{V}(\mathbf{f}_t) = T^{-2} \sum_{t,s} (\mathbf{f}_t - \bar{\mathbf{f}})(\mathbf{f}_s - \bar{\mathbf{f}})^\top$. By Davydov's inequality [1], for each $k \in \{1, \dots, K\}$ and $t, s \in \{1, \dots, T\}$, $|\mathbb{E}(f_{tk} f_{sk})^2| \lesssim \{\alpha(|t-s|)\}^{1/r_1} \{\mathbb{E}(|f_{tk}|^{2q_1})\}^{1/q_1} \{\mathbb{E}(|f_{sk}|^{2q_2})\}^{1/q_2}$, for some $q_1, q_2 > 0$ such that $1/r_1 + 1/q_1 + 1/q_2 = 1$, where $\alpha(\cdot)$ is the α -mixing coefficient. By Condition 7, $\mathbb{E}(|f_{tk}|^{q_1})$ and $\mathbb{E}(|f_{sk'}|^{q_2})$ exist for each $t \in \{1, \dots, T\}$ and $\alpha(|t-s|) < \exp(-C_1|t-s|^{r_1})$, so $|\mathbb{E}(f_{tk} f_{sk})^2| \lesssim \exp(-|t-s|)$. Thus, $\|\text{Cov}(\mathbf{f}_t, \mathbf{f}_s)\|_{\mathbb{F}} \lesssim \exp(-|t-s|)$ and $\|\mathcal{V}(\mathbf{f}_t) - \text{Var}(T^{-1} \sum_{t=1}^T \mathbf{f}_t)\|_{\mathbb{F}} = \|T^{-2} \sum_{t,s} \text{Cov}(\mathbf{f}_t, \mathbf{f}_s)\|_{\mathbb{F}} \lesssim T^{-2} \sum_{t=1}^T \exp(-t) \lesssim T^{-2}$. \square

Lemma 13. For each $i \in \{1, \dots, n\}$, with probability at least $1 - \delta$, $|\mathcal{V}(\hat{\mathbf{u}}_{it}) - \text{Var}(T^{-1} \sum_{t=1}^T u_{it})| \lesssim T^{-1} [\{\sqrt{n}T\}^{-1} + n^{-3/2+\alpha} + \{T^{1/2}n^{3/2}\}^{-1}p + \{\sqrt{n^2T}\}^{-1}\{(T+p^2)\ln(n)\}^{1/4} + n^{-1}J^{-\kappa/2}\}\{1 + \sqrt{\ln(21/\delta)}\}]$, where $\mathcal{V}(\hat{\mathbf{u}}_{it})$ is defined in Section 2.2.2.

Proof. Denote $\hat{\mathbf{U}} = \{\hat{\mathbf{u}}_{it}\}_{i=1, t=1}^{n,T}$. Note that $\mathbf{U} - \hat{\mathbf{U}} = (\hat{\mathbf{G}} - \mathbf{G}\mathbf{H}^{-1})(\hat{\mathbf{F}}^\top - \mathbf{H}\mathbf{F}^\top) + \mathbf{G}\mathbf{H}^{-1}(\hat{\mathbf{F}}^\top - \mathbf{H}\mathbf{F}^\top) + (\hat{\mathbf{G}} - \mathbf{G}\mathbf{H}^{-1})\mathbf{H}\mathbf{F}^\top$. By Propositions 1 and 3, with probability at least $1 - 20\delta$, $T^{-1}\|\hat{\mathbf{U}} - \mathbf{U}\|_{\mathbb{F}}^2 \lesssim \{n^{-1} + \{n^{1-2\alpha}T\}^{-1}p^2 + T^{-1}\sqrt{(T+p^2n^{2\alpha})\ln(n)} + J^{-\kappa}\}\{1 + \ln(1/\delta)\}^2$. Thus, similarly to the proof of Lemmas 11 and 12, with probability at least $1 - \delta$, $|\mathcal{V}(\hat{\mathbf{u}}_{it}) - \mathcal{V}(\mathbf{u}_{it})| \lesssim (nT)^{-1}[n^{-1/2} + (n^{1-2\alpha}T)^{-1/2}p + T^{-1/2}\{(T+p^2n^{2\alpha})\ln(n)\}^{1/4} + J^{-\kappa/2}\}\{1 + \sqrt{\ln(20/\delta)}\}]$ and $|\mathcal{V}(\mathbf{u}_{it}) - \text{Var}(T^{-1} \sum_{t=1}^T u_{it})| \lesssim \{\sqrt{nT^2}\}^{-1}\{1 + \sqrt{\ln(1/\delta)}\}$. The conclusion follows. \square

Lemma 14. With probability at least $1 - \delta$,

$$\|\mathbf{V}^{-1}(\hat{\mathbf{V}} - \mathbf{V})\|_2 \lesssim \left\{ \frac{\sqrt{J}}{n} + \frac{1}{\sqrt{n}} + \frac{1}{T} + \frac{p\sqrt{J}}{\sqrt{n^{1-2\alpha}T}} + \frac{1}{J^{(\kappa-1)/2}} \right\} \{1 + \ln(21/\delta)\}.$$

Proof. Recall that $\mathbf{V} = \mathbf{G}\text{Var}(T^{-1} \sum_{t=1}^T \mathbf{f}_t)\mathbf{G}^\top + \mathcal{D}$, so $\lambda_{\min}(\mathbf{V}) \geq \lambda_{\min}\{\mathbf{G}\text{Var}(T^{-1} \sum_{t=1}^T \mathbf{f}_t)\mathbf{G}^\top\} + \lambda_{\min}(\mathcal{D}) \gtrsim T^{-1}$. Note that $\hat{\mathbf{V}} - \mathbf{V} = \mathbf{G}\{\mathcal{V}(\hat{\mathbf{f}}_t) - \text{Var}(T^{-1} \sum_{t=1}^T \mathbf{f}_t)\}\mathbf{G}^\top + (\hat{\mathbf{G}} - \mathbf{G})\mathcal{V}(\hat{\mathbf{f}}_t)\hat{\mathbf{G}}^\top + \mathbf{G}\mathcal{V}(\hat{\mathbf{f}}_t)(\hat{\mathbf{G}} - \mathbf{G})^\top + (\hat{\mathcal{D}} - \mathcal{D})$. In addition, by the proof of Theorem 2 in Fan et al. [15], $\|\mathbf{G}\mathbf{V}^{-1}\mathbf{G}\|_2 = O(T)$. Thus,

$$\|\mathbf{V}^{-1}(\hat{\mathbf{V}} - \mathbf{V})\|_2 \leq \|\mathcal{V}(\hat{\mathbf{f}}_t) - \text{Var}(T^{-1} \sum_{t=1}^T \mathbf{f}_t)\|_{\mathbb{F}} + 2\|\mathcal{V}(\hat{\mathbf{f}}_t)\|_{\mathbb{F}}\|\hat{\mathbf{G}} - \mathbf{G}\|_{\mathbb{F}} + \|\hat{\mathcal{D}} - \mathcal{D}\|_{\mathbb{F}}.$$

From Lemmas 11 and 12, with probability at least $1 - \delta$,

$$\left\| \mathcal{V}(\hat{\mathbf{f}}_t) - \text{Var}\left(\frac{1}{T} \sum_{t=1}^T \mathbf{f}_t\right) \right\|_{\mathbb{F}} \lesssim \frac{1}{T} \left(\frac{1}{\sqrt{n}} + \frac{1}{T} + \frac{p}{\sqrt{n^{1-2\alpha}T}} + \frac{1}{J^{\kappa/2}} \right) \{1 + \sqrt{\ln(21/\delta)}\},$$

and

$$\|\hat{\mathcal{D}} - \mathcal{D}\|_{\mathbb{F}} \lesssim \frac{1}{T} \left[\frac{1}{T} + \frac{1}{n} + \frac{p}{n^{1-2\alpha}T^{1/2}} + \frac{\{(T+p^2n^{2\alpha})\ln(n)\}^{1/4}}{\sqrt{nT}} + \frac{1}{\sqrt{n}J^{\kappa/2}} \right] \{1 + \sqrt{\ln(21/\delta)}\}$$

which leads to the desired assertion by Lemma 13 and Theorem 1. \square

As a straightforward corollary to Lemma 14, with probability at least $1 - \delta$, $\|\hat{\mathbf{V}} - \mathbf{V}\|_{\mathbf{V}, \mathbb{F}} \lesssim \{n^{-1}\sqrt{J} + n^{-1/2} + T^{-1} + (\sqrt{n^{1-2\alpha}T})^{-1}p\sqrt{J} + J^{-(\kappa-1)/2}\}\sqrt{\ln(1/\delta)}$, where $\|\mathbf{A}\|_{\mathbf{S}, \mathbb{F}} := n^{-1/2}\|\mathbf{S}^{-1/2}\mathbf{A}\mathbf{S}^{-1/2}\|_{\mathbb{F}}$. If \mathbf{f}_t and \mathbf{u}_t are independent across t , then $\|\hat{\mathbf{V}} - \mathbf{V}\|_{\mathbf{V}, \mathbb{F}} \lesssim \{n^{-1}\sqrt{J} + p\sqrt{J}n^{-1/2+\alpha}T^{-1/2} + J^{-(\kappa-1)/2}\}\sqrt{\ln(1/\delta)}$, which mimics the optimal rate from Fan et al. [18] and Wang and Fan [43].

8.3. Discussions on legitimate preliminary $\hat{\boldsymbol{\beta}}^0$

In this section, we will discuss some preliminary estimators $\hat{\boldsymbol{\beta}}^0$ that satisfy the condition of TOPE, i.e., $\|\hat{\boldsymbol{\beta}}^0 - \boldsymbol{\beta}\|_2 = O_p(n^{-1/2+\alpha}T^{-1/2})$ for $\alpha \in [0, 1/2)$ in Section 2.2.2. In fact, Conditions 1, 2, and 6 guarantee the existence of such a preliminary $\hat{\boldsymbol{\beta}}^0$. We start with an OLS estimator based on an average version of (3) over time,

$$\hat{\boldsymbol{\beta}}^{\text{OLS}} = (\mathbb{Z}_0^\top \mathbb{Z}_0)^{-1} \mathbb{Z}_0^\top \bar{\mathbf{y}}. \quad (14)$$

Before showing that $\hat{\boldsymbol{\beta}}^{\text{OLS}}$ is a legitimate preliminary in Proposition 5, we first collect some technical results in Lemmas 15–17.

Lemma 15. Under Condition 7, $v_i(T) = T^{-1/2} \sum_{t=1}^T u_{it}$ is sub-exponential for each $i \in \{1, \dots, n\}$.

Proof. Note that $E(|u_{it}|^{4+\delta_1}) < \infty$ for any $t \in \{1, \dots, T\}$, $i \in \{1, \dots, n\}$ and $\delta_1 > 0$ and $\sum_{t=0}^{\infty} \alpha(T)^{1/3} < \sum_{t=0}^{\infty} \exp(-C T^{1/3}) < \infty$. By Theorem 4 in Tikhomirov [42], $|\Pr\{v_i(T) < s\} - \Pr\{W_i < s\}| \leq C_1 T^{-1/2} (1 + |s|)^{-4} \{\ln(T)\}^3$ for each $i \in \{1, \dots, n\}$ and any s , where $W \sim N(0, \sigma_i^2)$ and $\sigma_i^2 = E(u_{i1}^2) + 2 \sum_{t=2}^{\infty} E(u_{it} u_{i1})$. Thus, we have

$$\Pr\{|v_i(T)| > s\} = \Pr(|T^{-1/2} \sum_{t=1}^T u_{it}| > s) \leq 2 \exp\{-s^2/(2\sigma_i^2)\} + C_1 T^{-1/2} (1 + s)^{-4} \{\ln(T)\}^3$$

for any T and constants $C_1 > 0$. Furthermore, for any $k \in \{1, 2, \dots\}$,

$$\begin{aligned} E\{|v_i(T)|^k\} &= \int_0^1 \Pr\{|v_i(T)| > s^{1/k}\} ds + \int_1^{\infty} \Pr\{|v_i(T)| > s^{1/k}\} ds \\ &\leq 1 + 2(2\sigma_i^2)^{k/2} k \Gamma(k/2) + C_1 \pi T^{-1/2} \{\ln(T)\}^3 k!, \end{aligned}$$

so that $E\{\exp\{sv_i(T)\}\} \leq 1 + \sum_{k=2}^{\infty} |s|^k E\{|v_i(T)|^k\}/k! \lesssim \exp\{2\sigma_i^2 s^2 + C_1 \pi T^{-1/2} \{\ln(T)\}^3\}$ for $|s| < \min\{1/\sigma_i, 1\}$. The assertion follows from the definition of sub-exponential distributions. \square

Lemma 16. Under Conditions 1 and 7, for any $s > 0$, $p \times n$ matrix \mathbf{A} and $\sigma^2 = \max_i \sigma_i^2$ with σ_i^2 defined in Lemma 15, $\Pr\{\|\mathbf{A} \sum_{t=1}^T \mathbf{u}_t/T\|_2 > s \|\mathbf{A}\|_{\mathbb{F}}/\sqrt{T}\} < 2p \exp\{-s^2/(2\sigma^2)\}$.

Proof. Write $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_p)^T$, where $\mathbf{a}_1, \dots, \mathbf{a}_p$ are n -dimensional vectors. For each $m \in \{1, \dots, p\}$ and $w \geq 0$, by Conditions 1, 2, and 7, Lemma 15, and Corollary 4 in Samson [36], $\Pr(|\mathbf{a}_m^T \sum_{t=1}^T \mathbf{u}_t/T| \geq s) = \Pr(|\sum_{i=1}^n a_{mi} v_i(T)| \geq s\sqrt{T}) \leq 2 \exp\{-s^2 T/(2\sigma^2 \|\mathbf{a}_m\|_2^2)\}$. Hence, $\Pr\{|\mathbf{a}_m^T \sum_{t=1}^T \mathbf{u}_t/T| > \|\mathbf{a}_m\|_2 s/\sqrt{T}\} \leq 2 \exp\{-s^2/(2\sigma^2)\}$ for any $m \in \{1, \dots, p\}$ and $\Pr\{\|\mathbf{A} \sum_{t=1}^T \mathbf{u}_t/T\|_2 > \|\mathbf{A}\|_{\mathbb{F}} s/\sqrt{T}\} \leq 2p \exp\{-s^2/(2\sigma^2)\}$. \square

Conclusion in Lemma 16 remains valid for correlated $\{u_{it}\}_{t=1}^T$ over i . In fact, if one assumes cross-sectional dependence of $\{u_{it}\}$ over i by letting $\max_{j \leq n} \sum_{i=1}^n |E(u_{it} u_{jt})| < C_2$, $\max_{i \leq n} \sum_{k=1}^n \sum_{m=1}^n \sum_{t=1}^T \sum_{s=1}^T |\text{cov}(u_{it} u_{kt}, u_{is} u_{ms})| < C_2$, and $(nT)^{-1} \sum_{i=1}^n \sum_{j=1}^n \sum_{t=1}^T \sum_{s=1}^T |E(u_{it} u_{js})| < C_2$ for some $C_2 > 0$, Corollary 4 in Samson [36] still applies.

Lemma 17. For $p \times K$ matrix \mathbf{A} , under Conditions 1 and 7, $\Pr\{\|\mathbf{A} \sum_{t=1}^T \mathbf{f}_t/T\|_2 > s \|\mathbf{A}\|_{\mathbb{F}}/\sqrt{T}\} \leq 2pC_3 \exp(-C_4 s^2/2)$ for constants $C_3, C_4 > 0$.

Proof. The proof is similar to that of Lemma 16 and omitted here. \square

Proposition 5. Under Conditions 1, 2, and 6, with probability at least $1 - \delta$,

$$\|\hat{\boldsymbol{\beta}}^{\text{OLS}} - \boldsymbol{\beta}\|_2^2 \lesssim \frac{p^2}{n^{1-2\alpha} T} \ln(1/\delta).$$

Proof. Combining (14) and $\bar{\mathbf{y}} = \mathbb{Z}_0^T \boldsymbol{\beta} + \mathbf{G} T^{-1} \sum_{t=1}^T \mathbf{f}_t + T^{-1} \sum_{t=1}^T \mathbf{u}_t$, we have $\hat{\boldsymbol{\beta}}^{\text{OLS}} = \boldsymbol{\beta} + (\mathbb{Z}_0^T \mathbb{Z}_0)^{-1} \mathbb{Z}_0^T \mathbf{G} (T^{-1} \sum_{t=1}^T \mathbf{f}_t) + (\mathbb{Z}_0^T \mathbb{Z}_0)^{-1} \mathbb{Z}_0^T (T^{-1} \sum_{t=1}^T \mathbf{u}_t) \equiv \boldsymbol{\beta} + \text{(I)} + \text{(II)}$, where $\text{(I)} = (\mathbb{Z}_0^T \mathbb{Z}_0)^{-1} \mathbb{Z}_0^T \mathbf{G} (T^{-1} \sum_{t=1}^T \mathbf{f}_t)$ and $\text{(II)} = (\mathbb{Z}_0^T \mathbb{Z}_0)^{-1} \mathbb{Z}_0^T (T^{-1} \sum_{t=1}^T \mathbf{u}_t)$. By Condition 6, with probability 1, $\|\mathbf{P}_Z \mathbf{G}\|_{\mathbb{F}}^2 \lesssim n^{2\alpha}$. In addition, eigenvalues of $n^{-1} \mathbb{Z}_0^T \mathbb{Z}_0$ is bounded away from 0 and infinity almost surely by Condition 6(i). Thus, eigenvalues of $(n^{-1} \mathbb{Z}_0^T \mathbb{Z}_0)^{-1}$ are bounded away from 0 and infinity almost surely. That is, $\|(\mathbb{Z}_0^T \mathbb{Z}_0)^{-1} \mathbb{Z}_0^T\|_{\mathbb{F}}^2 = \text{tr}\{(\mathbb{Z}_0^T \mathbb{Z}_0)^{-1}\} \lesssim n^{-1} p$, and thus, $\|(\mathbb{Z}_0^T \mathbb{Z}_0)^{-1} \mathbb{Z}_0^T \mathbf{G}\|_{\mathbb{F}}^2 \lesssim \|(\mathbb{Z}_0^T \mathbb{Z}_0)^{-1} \mathbb{Z}_0^T\|_{\mathbb{F}}^2 \|\mathbf{P}_Z \mathbf{G}\|_{\mathbb{F}}^2 \lesssim n^{-1+2\alpha} p^2$ by Cauchy-Schwarz inequality. In light of Lemma 17, we have $\Pr\{\|(\text{I})\|_2 > s T^{-1/2} \|(\mathbb{Z}_0^T \mathbb{Z}_0)^{-1} \mathbb{Z}_0^T \mathbf{G}\|_{\mathbb{F}}\} < C_1 \exp(-C_2 s^2)$. By Lemma 16, it holds $\Pr\{\|(\text{II})\|_2 > s T^{-1/2} \|(\mathbb{Z}_0^T \mathbb{Z}_0)^{-1} \mathbb{Z}_0^T\|_{\mathbb{F}}\} < C_1 \exp(-C_2 s^2)$. Thus, we have $\Pr\{\|\hat{\boldsymbol{\beta}}^{\text{OLS}} - \boldsymbol{\beta}\|_2 > s T^{-1/2} \{\|(\mathbb{Z}_0^T \mathbb{Z}_0)^{-1} \mathbb{Z}_0^T \mathbf{G}\|_{\mathbb{F}} + \|(\mathbb{Z}_0^T \mathbb{Z}_0)^{-1} \mathbb{Z}_0^T\|_{\mathbb{F}}\} < 2C_1 \exp(-C_2 s^2)$. \square

Proposition 5 implies that $\hat{\boldsymbol{\beta}}^{\text{OLS}}$ is a legitimate preliminary for TOPE. Alternatively, one may consider the following choice on $\hat{\boldsymbol{\beta}}^0$. Rewrite $\mathbf{g}(\mathbf{x}_i)$ as $\mathbf{g}(\mathbf{x}_i) = \mathbf{A} \mathbf{z}_i + \mathbf{g}_0(\mathbf{x}_i)$, where \mathbf{A} is a $K \times p$ matrix and $\mathbf{z}_i = T^{-1} \sum_{t=1}^T \mathbf{z}_{it}$ is the average of \mathbf{z}_{it} over time. Then, model (1) can be rewritten as $y_{it} = \mathbf{z}_{it}^T \boldsymbol{\beta} + \mathbf{z}_i^T \boldsymbol{\eta}_t + \mathbf{g}_0(\mathbf{x}_i)^T \mathbf{f}_t + u_{it}$, where $\boldsymbol{\eta}_t = \mathbf{A}^T \mathbf{f}_t$. Under Condition 1, $\mathbf{g}_0(\mathbf{x}_i)^T \mathbf{f}_t + u_{it}$ is uncorrelated with the regressors \mathbf{z}_{it} . Hence, we can use the following random-effects GLS [37] to estimate

$(\beta, \eta_1, \dots, \eta_T)$ by $(\hat{\beta}^\top, \hat{\eta}_1^\top, \dots, \hat{\eta}_T^\top)^\top = (\mathbf{W}^\top \hat{\Sigma}_R^{-1} \mathbf{W})^{-1} \mathbf{W}^\top \hat{\Sigma}_R^{-1} \mathbf{y}$, where $\mathbf{y} = (y_{11}, \dots, y_{n1}, \dots, y_{1T}, \dots, y_{nT})^\top$,

$$\mathbf{W} = \begin{bmatrix} \mathbf{z}_{11}^\top & \mathbf{z}_{1\cdot}^\top & & \\ \vdots & \vdots & & \\ \mathbf{z}_{n1}^\top & \mathbf{z}_{n\cdot}^\top & & \\ \vdots & & \ddots & \\ \mathbf{z}_{1T}^\top & & & \mathbf{z}_{1\cdot}^\top \\ \vdots & & & \vdots \\ \mathbf{z}_{nT}^\top & & & \mathbf{z}_{n\cdot}^\top \end{bmatrix},$$

and $\hat{\Sigma}_R$ is an estimator of Σ_R , the covariance matrix of $\mathbf{v} = (\mathbf{g}_0(\mathbf{x}_1)^\top \mathbf{f}_1 + u_{11}, \dots, \mathbf{g}_0(\mathbf{x}_n)^\top \mathbf{f}_1 + u_{n1}, \dots, \mathbf{g}_0(\mathbf{x}_1)^\top \mathbf{f}_T + u_{1T}, \dots, \mathbf{g}_0(\mathbf{x}_n)^\top \mathbf{f}_T + u_{nT})^\top$. Under [Condition 1](#), Σ_R is a block diagonal matrix $\text{diag}(\Sigma_{R,1}, \dots, \Sigma_{R,T})$ with $\Sigma_{R,t} = E\{(\mathbf{g}_0(\mathbf{x}_1), \dots, \mathbf{g}_0(\mathbf{x}_n))^\top (\mathbf{g}_0(\mathbf{x}_1), \dots, \mathbf{g}_0(\mathbf{x}_n))\} + \sigma_u^2 \mathbf{I}_n$ for each $t \in \{1, \dots, T\}$, where $\text{var}(u_{it}) = \sigma_u^2$. There are a variety of estimators of $\Sigma_{R,1}$. For instance, Bai [3] and Schmidheiny and Basel [37] estimated $\hat{\Sigma}_{R,1}$ by first estimating \mathbf{v} , which is achieved via the OLS estimator. This is the so-called feasible GLS estimator [3,26,27] and can be extended to the iterative feasible GLS estimator [3,34]. That is, we can update $\hat{\Sigma}_{R,1}^{\text{new}}$ using $(\hat{\beta}^{\text{old}}, \hat{\eta}_1^{\text{old}}, \dots, \hat{\eta}_T^{\text{old}})$ from the previous step and iteratively update $(\hat{\beta}^{\text{new}}, \hat{\eta}_1^{\text{new}}, \dots, \hat{\eta}_T^{\text{new}})$ using the update $\hat{\Sigma}_{R,1}^{\text{new}}$. The update $(\hat{\beta}^{\text{new}}, \hat{\eta}_1^{\text{new}}, \dots, \hat{\eta}_T^{\text{new}})$ admits the following shrinkage of errors.

Proposition 6 (Lemma 1 in Phillips [34]). Under Conditions C1 to C3 in Phillips [34], if $T \geq p + 1$ and $\mathbf{A}_0 = E(\mathbf{W}^\top \Sigma_R^{-1} \mathbf{W})$ is nonsingular, $\sqrt{n}\{(\hat{\beta}^{\text{new}})^\top, (\hat{\eta}_1^{\text{new}})^\top, \dots, (\hat{\eta}_T^{\text{new}})^\top\} - (\beta^\top, \eta_1^\top, \dots, \eta_T^\top)^\top\} = 2(T-1)^{-1} \sqrt{nT}\{(\hat{\beta}^{\text{old}})^\top, (\hat{\eta}_1^{\text{old}})^\top, \dots, (\hat{\eta}_T^{\text{old}})^\top\} - (\beta^\top, \eta_1^\top, \dots, \eta_T^\top)^\top\} \psi \mathbf{A}_0^{-1} \psi + o_p(1)$, where ψ is given in Phillips [34].

Together along with $\|\hat{\beta} - \beta\|_2 \leq \|(\hat{\beta}^\top, \hat{\eta}_1^\top, \dots, \hat{\eta}_T^\top)^\top - (\beta^\top, \eta_1^\top, \dots, \eta_T^\top)^\top\|_2$, [Proposition 6](#) implies that the iterative feasible GLS estimator improves as the iteration grows. Thus, upon some iterations, the iterative feasible GLS estimator also provide a legitimate preliminary estimator for TOPE.

CRediT authorship contribution statement

Lyuou Zhang: Methodology, Theoretical analysis, Numerical study, Data analysis, Writing – original draft, Writing – review & editing. **Wen Zhou:** Conceptualization, Methodology, Theoretical analysis, Writing – original draft, Writing – review & editing, Supervision. **Haonan Wang:** Conceptualization, Methodology, Theoretical analysis, Writing – original draft, Writing – review & editing, Supervision.

Acknowledgments

The authors thank the Editor-in-Chief, an Associate Editor, and two reviewers for many helpful and constructive comments. The work of Wen Zhou was partially supported by Department of Energy, USA grant DE-SC0018344 and National Science Foundation, USA grants IIS-1545994 and IOS-1922701. The research of Haonan Wang was partially supported by National Science Foundation, USA grants DMS-1737795, DMS-1923142 and CNS-1932413.

Appendix A. Supplementary data

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.jmva.2021.104786>.

References

- [1] K.B. Athreya, S.N. Lahiri, Measure Theory and Probability Theory, Springer Science & Business Media, New York, 2006.
- [2] J. Bai, Inferential theory for factor models of large dimensions, *Econometrica* 71 (2003) 135–171.
- [3] J. Bai, Panel data models with interactive fixed effects, *Econometrica* 77 (2009) 1229–1279.
- [4] J. Bai, K. Li, Theory and methods of panel data models with interactive effects, *Ann. Statist.* 42 (2014) 142–170.
- [5] J. Bai, S. Ng, Principal components estimation and identification of static factors, *J. Econometrics* 176 (2013) 18–29.
- [6] B.H. Baltagi, *Econometrics*, fourth ed., Springer-Verlag, New York, 2008.
- [7] M. Barigozzi, H. Cho, P. Fryzlewicz, Simultaneous multiple change-point and factor analysis of high-dimensional time series, *J. Econometrics* 206 (2018) 87–225.
- [8] D. Bianchi, M. Billio, R. Casarin, M. Guidolin, Modeling systemic risk with Markov switching graphical SUR models, *J. Econometrics* 210 (2019) 58–74.
- [9] P. Billingsley, *Probability and Measure*, John Wiley & Sons, New York, 2012.
- [10] M. Cao, W. Zhou, F.J. Breidt, G. Peers, Large scale maximum average power multiple inference on time-course count data with application to RNA-seq analysis, *Biometrics* 76 (2020) 9–22.

- [11] G. Chamberlain, M. Rothschild, Arbitrage, factor structure, and mean-variance analysis on large asset markets, *Econometrica* 51 (1983) 1281–1304.
- [12] G. Connor, M. Hagmann, O. Linton, Efficient semiparametric estimation of the fama-french model and extensions, *Econometrica* 80 (2012) 713–754.
- [13] G. Connor, O. Linton, Semiparametric estimation of a characteristic-based factor model of common stock returns, *J. Empir. Financ.* 14 (2007) 694–717.
- [14] Z.J. Daye, J. Chen, H. Li, High-dimensional heteroscedastic regression with an application to eQTL data analysis, *Biometrics* 68 (2012) 316–326.
- [15] J. Fan, Y. Fan, J. Lv, High dimensional covariance matrix estimation using a factor model, *J. Econometrics* 147 (2008) 186–197.
- [16] J. Fan, T. Huang, Profile likelihood inference on semiparametric varying-coefficient partially linear models, *Bernoulli* 11 (2005) 1031–1057.
- [17] J. Fan, Y. Liao, M. Mincheva, High dimensional covariance matrix estimation in approximate factor models, *Ann. Statist.* 39 (2011) 3320–3356.
- [18] J. Fan, Y. Liao, M. Mincheva, Large covariance estimation by thresholding principal orthogonal complements, *J. R. Stat. Soc. Ser. B Stat. Methodol.* 75 (2013) 603–680.
- [19] J. Fan, Y. Liao, W. Wang, Projected principal component analysis in factor models, *Ann. Statist.* 44 (2016) 219–254.
- [20] M. Hallin, R. Liška, Dynamic factors in the presence of block structure, *J. Econometrics* 163 (2011) 29–41.
- [21] T. Hastie, R. Tibshirani, Generalized additive models, *Statist. Sci.* 1 (1986) 297–318.
- [22] R. Häuser, H. Liang, J.D. Meeker, H. Su, S.W. Thurston, Empirical likelihood based inference for additive partial linear measurement error models, *Stat. Interface* 2 (2009) 83–90.
- [23] J.Z. Huang, C.O. Wu, L. Zhou, Polynomial spline estimation and inference for varying coefficient models with longitudinal data, *Statist. Sinica* 14 (2004) 763–788.
- [24] J. Jiang, REML estimation: asymptotic behavior and related topics, *Ann. Statist.* 24 (1996) 255–286.
- [25] I.T. Jolliffe, *Principal Component Analysis*, second ed., Springer-Verlag, New York, 2002.
- [26] C. Lam, Q. Yao, Factor modeling for high-dimensional time series: inference for the number of factors, *Ann. Statist.* 40 (2012) 694–726.
- [27] J.T. Leek, J.D. Storey, Capturing heterogeneity in gene expression studies by surrogate variable analysis, *PLoS Genet.* 3 (2007) e161.
- [28] X. Liang, T. Zou, B. Guo, S. Li, H. Zhang, S. Zhang, H. Huang, S.X. Chen, Assessing Beijing's PM_{2.5} pollution: severity, weather impact, APEC and winter heating, *Proc. R. Soc. Lond. Ser. A Math. Phys. Eng. Sci.* 471 (2015) 2015.0257.
- [29] J. Lindström, A.A. Szpiro, P.D. Sampson, A.P. Oron, M. Richards, T.V. Larson, L. Sheppard, A flexible spatio-temporal model for air pollution with spatial and spatio-temporal covariates, *Environ. Ecol. Stat.* 21 (2014) 411–433.
- [30] G.G. Lorentz, *Approximation of Functions*, second ed., American Mathematical Society, 1986.
- [31] Z. Lu, D.J. Steinskog, D. Tjøstheim, Q. Yao, Adaptively varying-coefficient spatiotemporal models, *J. R. Stat. Soc. Ser. B* 71 (2009) 859–880.
- [32] F. Merlevède, M. Peligrad, E. Rio, A Bernstein type inequality and moderate deviations for weakly dependent sequences, *Probab. Theory Related Fields* 151 (2011) 435–474.
- [33] G. Motta, H. Ombao, Evolutionary factor analysis of replicated time series, *Biometrics* 68 (2012) 825–836.
- [34] R.F. Phillips, Iterated feasible generalized least-squares estimation of augmented dynamic panel data models, *J. Bus. Econom. Statist.* 28 (2010) 410–422.
- [35] P.M. Robinson, Root- N -consistent semiparametric regression, *Econometrica* 56 (1988) 931–954.
- [36] P.-M. Samson, Concentration of measure inequalities for Markov chains and Φ -mixing processes, *Ann. Probab.* 28 (2000) 416–461.
- [37] K. Schmidheiny, U. Basel, Panel data: fixed and random effects, *Short Guides To Microeconometrics* 7 (2011) 2–7.
- [38] J.H. Stock, M.W. Watson, Forecasting using principal components from a large number of predictors, *J. Amer. Statist. Assoc.* 97 (2002) 1167–1179.
- [39] C.J. Stone, Additive regression and other nonparametric models, *Ann. Statist.* 10 (1985) 689–705.
- [40] J.D. Storey, W. Xiao, J.T. Leek, R.G. Tompkins, R.W. Davis, Significance analysis of time course microarray experiments, *Proc. Natl. Acad. Sci. USA* 102 (2005) 12837–12842.
- [41] Z. Tan, G. Qin, H. Zhou, Estimation of a partially linear additive model for data from an outcome-dependent sampling design with a continuous outcome, *Biostatistics* 17 (2016) 663–676.
- [42] A.N. Tikhomirov, On the convergence rate in the central limit theorem for weakly dependent random variables, *Theory Probab. Appl.* 25 (1981) 790–809.
- [43] W. Wang, J. Fan, Asymptotics of empirical eigenstructure for high dimensional spiked covariance, *Ann. Statist.* 45 (2017) 1342–1374.
- [44] F. Wang, H. Wang, Modelling non-stationary multivariate time series of counts via common factors, *J. R. Stat. Soc. Ser. B Stat. Methodol.* 80 (2018) 769–791.
- [45] J. Wang, Q. Zhao, T. Hastie, A.B. Owen, Confounder adjustment in multiple hypothesis testing, *Ann. Statist.* 45 (2017) 1863–1894.
- [46] L. Wang, X.-H. Zhou, Assessing the adequacy of variance function in heteroscedastic regression models, *Biometrics* 63 (2007) 1218–1225.