

TENSOR FACTOR MODEL ESTIMATION BY ITERATIVE PROJECTION

BY YUEFENG HAN^{1,*}, RONG CHEN^{1,†}, DAN YANG², AND CUN-HUI ZHANG^{1,‡}

¹*Department of Statistics, Rutgers University, Piscataway, NJ 08854, USA, *yuefeng.han@rutgers.edu;
†rongchen@stat.rutgers.edu; ‡czhang@stat.rutgers.edu*

²*Faculty of Business and Economics, The University of Hong Kong, Hong Kong, dyanghku@hku.hk*

Tensor time series, which is a time series consisting of tensorial observations, has become ubiquitous. It typically exhibits high dimensionality. One approach for dimension reduction is to use a factor model structure, in a form similar to Tucker tensor decomposition, except that the time dimension is treated as a dynamic process with a time dependent structure. In this paper we introduce two approaches to estimate such a tensor factor model by using iterative orthogonal projections of the original tensor time series. These approaches extend the existing estimation procedures and improve the estimation accuracy and convergence rate significantly as proven in our theoretical investigation. Our algorithms are similar to the higher order orthogonal projection method for tensor decomposition, but with significant differences due to the need to unfold tensors in the iterations and the use of autocorrelation. Consequently, our analysis is significantly different from the existing ones. Computational and statistical lower bounds are derived to prove the optimality of the sample size requirement and convergence rate for the proposed methods. Simulation study is conducted to further illustrate the statistical properties of these estimators.

1. Introduction. Motivated by a diverse range of modern scientific applications, analysis of tensors, or multi-dimensional arrays, has emerged as one of the most important and active research areas in statistics, computer science, and machine learning. Large tensors are encountered in genomics (Alter and Golub, 2005, Omberg, Golub and Alter, 2007), neuroimaging analysis (Zhou, Li and Zhu, 2013, Sun and Li, 2017), recommender systems (Bi, Qu and Shen, 2018), computer vision (Liu et al., 2012), community detection (Anandkumar et al., 2014), among others. High-order tensors often bring about high dimensionality and impose significant computational challenges. For example, functional MRI produces a time series of 3-dimensional brain images, typically consisting of hundreds of thousands of voxels observed over time. Previous work has developed various tensor-based methods for independent and identically distributed (i.i.d.) tensor data or tensor data with i.i.d. noise. However, as far as we know, the statistical framework for general tensor time series data was not well studied in the literature.

Factor analysis is one of the most useful tools for understanding common dependence among multi-dimensional outputs. Over the past decades, vector factor models have been extensively studied in the statistics and economics communities. For instance, Chamberlain and Rothschild (1983), Bai and Ng (2002), Stock and Watson (2002) and Bai (2003) developed the static factor model using principal component analysis (PCA). They assumed that the common factors must have impact on most of the time series, and weak serial dependence is allowed for the idiosyncratic noise process. Fan, Liao and Mincheva (2011, 2013), Fan, Liu and Wang (2018) established large covariance matrix estimation based on the static

MSC2020 subject classifications: Primary 62H25, 62H12; secondary 62R07.

Keywords and phrases: high-dimensional tensor data, factor model, orthogonal projection, time series, Tucker decomposition.

factor model. The static factor model has been further extended to the dynamic factor model in [Forni et al. \(2000\)](#). The latent factors are assumed to follow a time series process, which is commonly taken to be a vector autoregressive process. [Fan, Liao and Wang \(2016\)](#) studied semi-parametric factor models through projected principal component analysis. [Pena and Box \(1987\)](#), [Pan and Yao \(2008\)](#), [Lam, Yao and Bathia \(2011\)](#) and [Lam and Yao \(2012\)](#) adopted another type of factor model. They assumed that the latent factors capture all dynamics of the observed process, and thus the idiosyncratic noise process has no serial dependence. We will adopt this approach. We note that the factor process may have complex dynamic behavior, resulting in complex dynamics of the observed tensor, even with white additive noise process. Of course, when all the dynamics of the observed tensor process are ‘forced’ to be included in the signal process induced by the factor process, situations may arise in which some factors are ‘weak’ (or have impact on a small portion of the observed series in the tensor). This leads us to consider the ‘signal strength’ in our investigation.

Although there have been significant efforts in developing methodologies and theories for vector factor models, there is a paucity of literature on matrix- or tensor-valued time series. [Wang, Liu and Chen \(2019\)](#) proposed a matrix factor model for matrix-valued time series, which explores the matrix structure. [Chen, Tsay and Chen \(2019\)](#) established a general framework for incorporating domain and prior knowledge in the matrix factor model through linear constraints. [Chen and Chen \(2019\)](#) applied the matrix factor model to the dynamic transport network. [Chen, Fan and Li \(2020\)](#) developed an inferential theory of the matrix factor model under a different setting from that in [Wang, Liu and Chen \(2019\)](#).

Recently, [Chen, Yang and Zhang \(2019\)](#) introduced a factor approach for analyzing high dimensional dynamic tensor time series in the form

$$(1.1) \quad \mathcal{X}_t = \mathcal{M}_t + \mathcal{E}_t,$$

where $\mathcal{X}_1, \dots, \mathcal{X}_T \in \mathbb{R}^{d_1 \times \dots \times d_K}$ are the observed tensor time series, \mathcal{M}_t and \mathcal{E}_t are the corresponding signal and noise components of \mathcal{X}_t , respectively. The goal is to estimate the unknown signal tensor \mathcal{M}_t from the tensor time series data. Following [Lam and Yao \(2012\)](#), it is assumed that the signal tensor accommodates all dynamics, making the idiosyncratic noise \mathcal{E}_t uncorrelated (white) across time. It is further assumed that \mathcal{M}_t lives in a lower dimensional space and has certain multilinear decomposition. Specifically, we assume that \mathcal{M}_t satisfies a Tucker-type decomposition and model (1.1) can be written as

$$(1.2) \quad \mathcal{X}_t = \mathcal{F}_t \times_1 A_1 \times_2 \dots \times_K A_K + \mathcal{E}_t,$$

where A_k is the deterministic loading matrix of size $d_k \times r_k$ and $r_k \ll d_k$, and the core tensor \mathcal{F}_t itself is a latent tensor factor process of dimension $r_1 \times \dots \times r_K$. Here the k -mode product of $\mathcal{X} \in \mathbb{R}^{d_1 \times d_2 \times \dots \times d_K}$ with a matrix $U \in \mathbb{R}^{d_k \times d_k}$, denoted as $\mathcal{X} \times_k U$, is an order K -tensor of size $d_1 \times \dots \times d_{k-1} \times d_k' \times d_{k+1} \times \dots \times d_K$ such that

$$(\mathcal{X} \times_k U)_{i_1, \dots, i_{k-1}, j, i_{k+1}, \dots, i_K} = \sum_{i_k=1}^{d_k} \mathcal{X}_{i_1, i_2, \dots, i_K} U_{j, i_k}.$$

The core tensor \mathcal{F}_t is usually much smaller than \mathcal{X}_t in dimension. This structure provides an effective dimension reduction, as all the comovements of individual time series in \mathcal{X}_t are driven by \mathcal{F}_t . Without loss of generality, assume that A_k is of rank $r_k \ll d_k$. It should be noted that vector and matrix factor models can be viewed as special cases of our model since a vector time series is a tensor time series composed of a single fiber ($K = 1$), and a matrix time series is one composed of a single slice ($K = 2$).

[Chen, Yang and Zhang \(2019\)](#) proposed two estimation procedures, namely TOPUP and TIPUP, for estimating the column space spanned by the loading matrix A_k , for $k = 1, \dots, K$.

The two procedures are based on different auto-cross-product operations of the observed tensors \mathcal{X}_t to accumulate information, but they both utilize the assumption that the noise \mathcal{E}_t and \mathcal{E}_{t-h} , $h > 0$ are uncorrelated. The convergence rates of their estimators critically depend on $d = d_1 d_2 \dots d_K$, a potentially very large number as d_k , $k = 1, \dots, K$, are large. Often a large T , the length of the time series, is required for accurate estimation of the loading spaces.

In this paper we propose extensions of the TOPUP and TIPUP procedures, motivated by the following observation. Suppose that the loading matrices A_k are orthonormal with $A_k^\top A_k = I$, and we are given A_2, \dots, A_K . Let

$$\mathcal{Z}_t = \mathcal{X}_t \times_2 A_2^\top \times_3 \dots \times_K A_K^\top; \text{ and } \mathcal{E}_t^* = \mathcal{E}_t \times_2 A_2^\top \times_3 \dots \times_K A_K^\top;$$

Then (1.2) leads to

$$(1.3) \quad \mathcal{Z}_t = \mathcal{F}_t \times_1 A_1 + \mathcal{E}_t^*$$

where \mathcal{Z}_t is a $d_1 \times r_2 \times \dots \times r_K$ tensor. Since $r_k \ll d_k$, \mathcal{Z}_t is a much smaller tensor than \mathcal{X}_t . Under proper conditions on the combined noise tensor \mathcal{E}_t^* , the estimation of the loading space of A_1 based on \mathcal{Z}_t can be made significantly more accurate, as the convergence rate now depends on $d_1 r_2 \dots r_K$ rather than $d_1 d_2 \dots d_K$.

Of course, in practice we do not know A_2, \dots, A_K . Similar to backfitting algorithms, we propose an iterative algorithm. With a proper initial value, we iteratively estimate the loading space of A_k at iteration j based on

$$\mathcal{Z}_{t,k}^{(j)} = \mathcal{X}_t \times_1 \hat{A}_1^{(j)\top} \times_2 \dots \times_{k-1} \hat{A}_{k-1}^{(j)\top} \times_{k+1} \hat{A}_{k+1}^{(j-1)\top} \times_{k+2} \dots \times_K \hat{A}_K^{(j-1)\top},$$

using the estimate $\hat{A}_{k'}^{(j-1)}$, $k < k' \leq K$ obtained in the previous iteration and the estimate $\hat{A}_{k'}^{(j)}$, $1 \leq k' < k$, obtained in the current iteration. Our theoretical investigation shows that the iterative procedures for estimating A_1 can achieve the convergence rate as if all A_2, \dots, A_K are known and we indeed observe \mathcal{Z}_t that follows model (1.3). We call the procedure iTOPUP and iTIPUP, based on the matrix unfolding mechanism used, corresponding to TOPUP and TIPUP procedures. To be more specific, our algorithms have two steps: (i) We first use the estimated column space of factor loading matrices of TOPUP (resp. TIPUP) to construct the initial estimate of factor loading spaces; (ii) We then iteratively perform matrix unfolding of the auto-cross-moments of much smaller tensors $\mathcal{Z}_{t,k}^{(j)}$ to obtain the final estimators.

We note that the iterative procedure is related to higher order orthogonal iteration (HOOI) that has been widely studied in the literature; see, e.g., De Lathauwer, De Moor and Vandewalle (2000), Sheehan and Saad (2007), Liu et al. (2014), Zhang and Xia (2018), among others. However, most of the existing works are not designed for tensor time series. They do not consider the special role of the time mode nor the covariance structure in the time direction. Typically HOOI treats the signal part as fixed or deterministic. In this paper we treat the signal as dynamic in the sense that the core tensor \mathcal{F}_t in (1.2) is dynamic and the relationship between \mathcal{F}_t and the lagged \mathcal{F}_{t-h} is of interest. Our setting requires special treatment although each iteration of our iterative procedures also consists of power up and orthogonal projection operations. While HOOI applies the SVD directly to the matrix unfolding of the iteratively projected data, in our approach the SVD is applied to the matrix unfolding of the outer- and inner-auto-cross-product of the iteratively projected data, respectively in iTOPUP and iTIPUP. Although the iTOPUP algorithm proposed here can be reformulated as a twist of HOOI on the auto-cross-moment tensor, the iTIPUP algorithm is different and cannot be recast equivalently as HOOI. More importantly, the theoretical analysis and theoretical properties of the estimators are fundamentally different from those of HOOI, due to the dynamic structure of tensor time series and the need to use the auto-cross-product operation between

the SVD and data projection in each iteration. Different concentration inequalities are derived to study the performance bounds.

In this paper, we establish upper bounds on the estimation errors for both the iTOPUP and the iTIPUP, which are much sharper than the respective theoretical guarantees for TOPUP and TIPUP, demonstrating the benefits of using iterative projection. It is also shown that the number of iterations needed for convergence is of order no greater than $\log(d)$. We mainly focus on the cases where the tensor dimensions are large and of similar order. We also cover the cases where the ranks of the tensor factor process increase with the dimensions of the tensor time series.

Chen, Yang and Zhang (2019) showed that the TIPUP has a faster convergence rate in estimation error than the TOPUP, under a mild condition on the level of signal cancellation. In contrast, the theoretically guaranteed rate of convergence for the iTOPUP in this paper is of the same order or even faster than that for the iTIPUP under certain regularity conditions. Our results also suggest an interesting phenomenon. Using the iterative procedures, we find that the increase in either dimension or sample size can improve the estimation of the factor loading space of the tensor factor model with the tensor order $K \geq 2$. We believe that such a super convergence rate is new in the literature. Specifically, under proper regularity conditions, the convergence rate of the iterative procedures for estimating the space of A_k is $O_{\mathbb{P}}(T^{-1/2}d_{-k}^{-1/2})$, where $d_{-k} = \prod_{j \neq k} d_j$, while the existing rate for non-iterative procedures is $O_{\mathbb{P}}(T^{-1/2})$ for the vector factor model (Lam, Yao and Bathia, 2011) and the matrix/tensor factor models (Wang, Liu and Chen, 2019, Chen, Yang and Zhang, 2019). While the increase in the dimensions d_k ($k = 1, \dots, K$) does not improve the performance of the non-iterative estimators, it significantly improves that of the proposed iterative estimators.

In addition, we establish the computational lower bound for the estimation of the loading spaces of tensor factor models under the hardness assumption of certain instances of hyper-graphic planted clique detection problem. It shows that the sample size requirement (or signal to noise ratio condition) needed for using the TIPUP estimate as the initial values for the iterative procedures is unavoidable for any computationally manageable estimation procedure to achieve consistency, although the iterative procedures have faster convergence rates. Moreover, we provide a statistical lower bound which matches the rates of convergence of our iterative procedures under proper conditions.

Related work. We close this section by highlighting several recent papers on related topics. First, we draw attention to the work of Foster (1996), Fan, Liao and Wang (2016) and Chen et al. (2020). Chen et al. (2020) adopts a spectral initialization plus an iterative refinement step estimating procedure, so that our methods are related to theirs. However, due to the differences in problem setting and model assumptions, their estimation procedures, performance bounds and analytic techniques are all significantly different from ours. Foster (1996), Fan, Liao and Wang (2016) use the projection to the space spanned by the sieve bases without iteration. Rogers, Li and Russell (2013) assumes the tensor factor model in (1.2), with an additional specific AR structure on the dynamic of the factor process. The additional model structure led to an EM type of estimation approach, quite different from the approach we develop here. Wang, Zheng and Li (2021) concerns low rank tensor AR model and uses a nuclear norm penalty to enforce the low rank structure and optimization algorithms for estimation, again quite different from our approach.

The paper is organized as follows. Section 2.1 introduces basic notation and preliminaries of tensor analysis. We present the tensor factor model and the iTOPUP and iTIPUP procedures in Sections 2.2 and 2.3. Theoretical properties of the iTOPUP and iTIPUP are investigated in Section 3. Section 4 provides a brief summary. Numerical comparison of our iterative procedures and other methods, and all technical details are relegated to the Supplementary Material.

2. Tensor Factor Model by Orthogonal Iteration.

2.1. *Notation and preliminaries for tensor analysis.* Throughout this paper, for a vector $x = (x_1, \dots, x_p)^\top$, define $\|x\|_q = (x_1^q + \dots + x_p^q)^{1/q}$, $q \geq 1$. For a matrix $A = (a_{ij}) \in \mathbb{R}^{m \times n}$, write the SVD as $A = U\Sigma V^\top$, where $\Sigma = \text{diag}(\sigma_1(A), \sigma_2(A), \dots, \sigma_{\min\{m,n\}}(A))$, with the singular values $\sigma_1(A) \geq \sigma_2(A) \geq \dots \geq \sigma_{\min\{m,n\}}(A) \geq 0$ in descending order. The matrix spectral norm is denoted as $\|A\|_S = \sigma_1(A)$. Let $\sigma_{\min}(A)$ (resp. $\sigma_{\max}(A)$) be the smallest (resp. largest) nontrivial singular value of A . For two sequences of real numbers $\{a_n\}$ and $\{b_n\}$, write $a_n = O(b_n)$ (resp. $a_n \asymp b_n$) if there exists a constant C such that $|a_n| \leq C|b_n|$ (resp. $1/C \leq a_n/b_n \leq C$) for all sufficiently large n , and write $a_n = o(b_n)$ if $\lim_{n \rightarrow \infty} a_n/b_n = 0$. Write $a_n \lesssim b_n$ (resp. $a_n \gtrsim b_n$) if there exist a constant C such that $a_n \leq Cb_n$ (resp. $a_n \geq Cb_n$). Denote $a \wedge b = \min\{a, b\}$ and $a \vee b = \max\{a, b\}$. We use C, C_1, c, c_1, \dots to denote generic constants, whose actual values may vary from line to line.

For any two $m \times r$ matrices with orthonormal columns, say, U and \hat{U} , suppose the singular values of $U^\top \hat{U}$ are $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r \geq 0$. A natural measure of distance between the column spaces of U and \hat{U} is then

$$(2.1) \quad \|\hat{U}\hat{U}^\top - UU^\top\|_S = \sqrt{1 - \sigma_r^2},$$

which equals to the sine of the largest principle angle between the column spaces of U and \hat{U} .

For any two matrices $A \in \mathbb{R}^{m_1 \times r_1}, B \in \mathbb{R}^{m_2 \times r_2}$, denote the Kronecker product \odot as $A \odot B \in \mathbb{R}^{m_1 m_2 \times r_1 r_2}$. For any two tensors $\mathcal{A} \in \mathbb{R}^{m_1 \times m_2 \times \dots \times m_K}, \mathcal{B} \in \mathbb{R}^{r_1 \times r_2 \times \dots \times r_N}$, denote the tensor product \otimes as $\mathcal{A} \otimes \mathcal{B} \in \mathbb{R}^{m_1 \times \dots \times m_K \times r_1 \times \dots \times r_N}$, such that

$$(\mathcal{A} \otimes \mathcal{B})_{i_1, \dots, i_K, j_1, \dots, j_N} = (\mathcal{A})_{i_1, \dots, i_K} (\mathcal{B})_{j_1, \dots, j_N}.$$

Let $\text{vec}(\cdot)$ be the vectorization of matrices and tensors. The mode- k unfolding (or matricization) is defined as $\text{mat}_k(\mathcal{A})$, which maps a tensor \mathcal{A} to a matrix $\text{mat}_k(\mathcal{A}) \in \mathbb{R}^{m_k \times m_{-k}}$ where $m_{-k} = \prod_{j \neq k}^K m_j$. For example, if $\mathcal{A} \in \mathbb{R}^{m_1 \times m_2 \times m_3}$, then

$$(\text{mat}_1(\mathcal{A}))_{i, (j+m_2(k-1))} = (\text{mat}_2(\mathcal{A}))_{j, (k+m_3(i-1))} = (\text{mat}_3(\mathcal{A}))_{k, (i+m_1(j-1))} = \mathcal{A}_{ijk}.$$

For tensor $\mathcal{A} \in \mathbb{R}^{m_1 \times m_2 \times \dots \times m_K}$, the Hilbert Schmidt norm is defined as

$$\|\mathcal{A}\|_{\text{HS}} = \sqrt{\sum_{i_1=1}^{m_1} \dots \sum_{i_K=1}^{m_K} (\mathcal{A})_{i_1, \dots, i_K}^2}.$$

For a matrix, the Hilbert Schmidt norm is just the Frobenius norm. Define the tensor operator norm for an order-4 tensor $\mathcal{A} \in \mathbb{R}^{m_1 \times m_2 \times m_3 \times m_4}$,

$$\|\mathcal{A}\|_{\text{op}} = \max \left\{ \sum_{i_1, i_2, i_3, i_4} u_{i_1, i_2} \cdot u_{i_3, i_4} \cdot (\mathcal{A})_{i_1, i_2, i_3, i_4} : \|U_1\|_{\text{HS}} = \|U_2\|_{\text{HS}} = 1 \right\},$$

where $U_1 = (u_{i_1, i_2}) \in \mathbb{R}^{m_1 \times m_2}$ and $U_2 = (u_{i_3, i_4}) \in \mathbb{R}^{m_3 \times m_4}$.

2.2. *Tensor factor model.* Again, we consider as in (1.2)

$$\mathcal{X}_t = \mathcal{F}_t \times_1 A_1 \times_2 \dots \times_K A_K + \mathcal{E}_t.$$

Without loss of generality, assume that A_k is of rank r_k . A_k is not necessarily orthonormal, which is different from the classical Tucker decomposition (Tucker, 1966). Model (1.2) is unchanged if we replace $(A_1, \dots, A_K, \mathcal{F}_t)$ by $(A_1 H_1, \dots, A_K H_K, \mathcal{F}_t \times_{k=1}^K H_k^{-1})$ for any invertible $r_k \times r_k$ matrix H_k . Although $(A_1, \dots, A_K, \mathcal{F}_t)$ are not uniquely determined, the factor

loading space, that is, the linear space spanned by the columns of A_k , is uniquely defined. Denote the orthogonal projection to the column space of A_k as

$$(2.2) \quad P_k = P_{A_k} = A_k(A_k^\top A_k)^{-1}A_k^\top = U_k U_k^\top,$$

where U_k is the left singular matrix in the SVD $A_k = U_k \Lambda_k V_k^\top$. We use P_k to represent the factor loading space of A_k . Thus, our objective is to estimate P_k .

The canonical representation of the tensor times series (1.2) is written as

$$\mathcal{X}_t = \mathcal{F}_t^{(\text{cano})} \times_{k=1}^K U_k + \mathcal{E}_t,$$

where the diagonal and right singular matrices of A_k are absorbed into the canonical core tensor $\mathcal{F}_t^{(\text{cano})} = \mathcal{F}_t \times_{k=1}^K (\Lambda_k V_k^\top)$. In this canonical form, the loading matrices U_k are identifiable up to a rotation in general and up to a permutation and sign changes of the columns of U_k when the singular values are all distinct in the population version of the TOPUP or TIPUP methods, as we describe in Section 2.3 below. In what follows, we may identify the tensor time series in its canonical form, i.e. $A_k = U_k$, without explicit declaration.

We do not impose any specific structure for the dynamics of the core tensor factor process $\mathcal{F}_t \in \mathbb{R}^{r_1 \times \dots \times r_K}$ beyond the independence between the core process and the noise process, and we do not require any additional structure on the correlation among different time series fibers of the noise process \mathcal{E}_t . Because of this generality, our estimator is based on the tensor version of the lagged sample cross product $\hat{\Sigma}_h$, $h = 1, \dots, h_0$, where

$$(2.3) \quad \hat{\Sigma}_h = \hat{\Sigma}_h(\mathcal{X}_{1:T}) = \sum_{t=h+1}^T \frac{\mathcal{X}_{t-h} \otimes \mathcal{X}_t}{T-h} \in \mathbb{R}^{d_1 \times \dots \times d_K \times d_1 \times \dots \times d_K},$$

which is an order- $2K$ tensor. The population version of this tensor autocovariance is

$$\Sigma_h = \mathbb{E} \left(\sum_{t=h+1}^T \frac{\mathcal{X}_{t-h} \otimes \mathcal{X}_t}{T-h} \right) = \mathbb{E} \left(\sum_{t=h+1}^T \frac{\mathcal{M}_{t-h} \otimes \mathcal{M}_t}{T-h} \right).$$

Because $\mathcal{M}_t = \mathcal{M}_t \times_{k=1}^K P_k$ for all t ,

$$\Sigma_h = \Sigma_h \times_{k=1}^{2K} P_k = \mathbb{E} \left(\sum_{t=h+1}^T \frac{\mathcal{F}_{t-h} \otimes \mathcal{F}_t}{T-h} \right) \times_{k=1}^{2K} P_k A_k,$$

with the notation $A_k = A_{k-K}$ and $P_k = P_{k-K}$ for all $k > K$.

2.3. Estimating procedures. In this paper, we consider iterative estimation procedures to achieve sharper convergence rates than the TOPUP and TIPUP procedures proposed in [Chen, Yang and Zhang \(2019\)](#). We start with a quick description of their procedures as they serve as the starting point of our proposed iTOPUP and iTIPUP procedures. Note that the procedure in [Chen and Chen \(2019\)](#) and [Wang, Liu and Chen \(2019\)](#) is the non-iterative TOPUP.

(i). Time series Outer-Product Unfolding Procedure (TOPUP):

Let $\hat{\Sigma}_h$ be the sample autocovariance of the data $\mathcal{X}_{1:T} = (\mathcal{X}_1, \dots, \mathcal{X}_T)$ as in (2.3). Define

$$(2.4) \quad \text{TOPUP}_k = \left(\text{mat}_k(\hat{\Sigma}_h), h = 1, \dots, h_0 \right),$$

as a $d_k \times (d_{-k} h_0)$ matrix, where $d = \prod_{k=1}^K d_k$, $d_{-k} = d/d_k$ and h_0 is a predetermined positive integer. Here we note that TOPUP_k is a function mapping a tensor time series to a matrix. In TOPUP_k , the information from different time lags is accumulated, which is useful especially when the sample size T is small. A relatively small h_0 is typically used, since the autocorrelation is often at its strongest with small time lags. See Remark 3.8.

The TOPUP method performs SVD of (2.4) to obtain the truncated left singular matrices

$$\hat{U}_{k,m}^{TOPUP}(\mathcal{X}_{1:T}) = \text{LSVD}_m \left(\text{mat}_k(\hat{\Sigma}_h(\mathcal{X}_{1:T})), h = 1, \dots, h_0 \right),$$

where LSVD_m stands for the left singular matrix composed of the first m left singular vectors corresponding to the largest m singular values. Here we emphasize that $\hat{U}_{k,m}^{TOPUP}(\cdot)$ is treated as an operator that maps a tensor data set to a matrix of m columns. It will be applied to different transformed data sets. On the other hand, TOPUP_k is treated as fixed, based on the given $\mathcal{X}_{1:T}$ under study. Note that LSVD can be obtained using eigen decomposition as well. For simplicity, we write

$$(2.5) \quad \text{UTOPUP}_k(\mathcal{X}_{1:T}, r_k) = \hat{U}_{k,r_k}^{TOPUP}(\mathcal{X}_{1:T}),$$

where r_k is the mode- k rank. Again, we emphasize that $\text{UTOPUP}_k(\cdot)$ takes input of a tensor time series of length T with the target mode- k having dimension d_k and rank r_k , and produces an output matrix of size $d_k \times r_k$ as the estimate of the mode- k loading matrix U_k .

By (1.2) and (2.3), the expectation of (2.4) satisfies

$$(2.6) \quad \mathbb{E}[\text{TOPUP}_k] \\ = A_k \text{mat}_k \left(\sum_{t=h+1}^T \mathbb{E} \left(\frac{\mathcal{F}_{t-h} \otimes \mathcal{F}_t}{T-h} \right) \times_{l=1}^{k-1} A_l \times_{l=k+1}^{2K} A_l, h = 1, \dots, h_0 \right),$$

so that the TOPUP is expected to be consistent in estimating the column space of A_k .

(ii). Time series Inner-Product Unfolding Procedure (TIPUP):

Similar to (2.4), define a $d_k \times (d_k h_0)$ matrix as

$$(2.7) \quad \text{TIPUP}_k = \left(\sum_{t=h+1}^T \frac{\text{mat}_k(\mathcal{X}_{t-h}) \text{mat}_k^\top(\mathcal{X}_t)}{T-h}, h = 1, \dots, h_0 \right),$$

which replaces the tensor product by the inner product through (2.3) in (2.4). The TIPUP method performs SVD:

$$\hat{U}_{k,m}^{TIPUP}(\mathcal{X}_{1:T}) = \text{LSVD}_m \left(\sum_{t=h+1}^T \frac{\text{mat}_k(\mathcal{X}_{t-h}) \text{mat}_k^\top(\mathcal{X}_t)}{T-h}, h = 1, \dots, h_0 \right),$$

for $k = 1, \dots, K$. Again, $\hat{U}_{k,m}^{TIPUP}(\cdot)$ is treated as an operator. For simplicity, we write

$$(2.8) \quad \text{UTIPUP}_k(\mathcal{X}_{1:T}, r_k) = \hat{U}_{k,r_k}^{TIPUP}(\mathcal{X}_{1:T}).$$

where r_k is the mode- k rank. Note that

$$(2.9) \quad \mathbb{E}[\text{TIPUP}_k] \\ = \text{mat}_1 \left(\langle \sum_h \mathcal{I}_{k,k+K} \rangle_{\{k,k+K\}^c}, h = 1, \dots, h_0 \right) \\ = A_k \text{mat}_1 \left(\left\langle \sum_{t=h+1}^T \mathbb{E} \left(\frac{\mathcal{F}_{t-h} \otimes \mathcal{F}_t}{T-h} \right) \times_{l \neq k, 1 \leq l \leq 2K} A_l, \mathcal{I}_{k,k+K} \right\rangle_{\{k,k+K\}^c}, h = 1, \dots, h_0 \right),$$

where $\mathcal{I}_{k,k+K}$ is an order- $2K$ tensor with elements $(\mathcal{I}_{k,k+K})_{\mathbf{i}, \mathbf{j}} = I\{\mathbf{i}_{-k} = \mathbf{j}_{-k}\}$, $\mathbf{i} = (i_1, \dots, i_K)$, $\mathbf{j} = (j_1, \dots, j_K)$, $\mathbf{i}_{-k} = (i_1, \dots, i_{k-1}, i_{k+1}, \dots, i_K)$, $\mathbf{j}_{-k} = (j_1, \dots, j_{k-1}, j_{k+1}, \dots, j_K)$.

Here, $\langle \cdot, \cdot \rangle_{\{k, k+K\}^c}$ is defined as an inner product summation over all indices other than $\{k, k+K\}$.

(iii). **iTOPUP and iTIPUP:** Next we describe a generic iterative procedure under the motivation described in Section 1. Its pseudo-code is provided in Algorithm 1. It incorporates two estimators/operators UINIT and UITER that map a tensor time series to an estimate of the loading matrix. The UTOPUP and UTIPUP operators in (2.5) and (2.8) are examples of such operators.

Algorithm 1 A generic iterative algorithm

- 1: Input: $\mathcal{X}_t \in \mathbb{R}^{d_1 \times \dots \times d_K}$ for $t = 1, \dots, T$, r_k for all $k = 1, \dots, K$, the tolerance parameter $\epsilon > 0$, the maximum number of iterations J , and the UINIT and UITER operators.
- 2: Let $j = 0$, initiate via applying UINIT on $\{\mathcal{X}_{1:T}\}$, for $k = 1, \dots, K$, to obtain

$$\hat{U}_k^{(0)} = \text{UINIT}_k(\mathcal{X}_{1:T}, r_k).$$

- 3: **repeat**

- 4: Let $j = j + 1$. At the j -th iteration, for $k = 1, \dots, K$, given previous estimates $(\hat{U}_{k+1}^{(j-1)}, \dots, \hat{U}_K^{(j-1)})$ and $(\hat{U}_1^{(j)}, \dots, \hat{U}_{k-1}^{(j)})$, sequentially calculate,

$$\mathcal{Z}_{t,k}^{(j)} = \mathcal{X}_t \times_1 (\hat{U}_1^{(j)})^\top \times_2 \dots \times_{k-1} (\hat{U}_{k-1}^{(j)})^\top \times_{k+1} (\hat{U}_{k+1}^{(j-1)})^\top \times_{k+2} \dots \times_K (\hat{U}_K^{(j-1)})^\top,$$

for $t = 1, \dots, T$. Perform UITER on the new tensor time series $\mathcal{Z}_{1:T,k}^{(j)} = (\mathcal{Z}_{1,k}^{(j)}, \dots, \mathcal{Z}_{T,k}^{(j)})$.

$$\hat{U}_k^{(j)} = \text{UITER}_k(\mathcal{Z}_{1:T,k}^{(j)}, r_k).$$

- 5: **until** $j = J$ or

$$\max_{1 \leq k \leq K} \|\hat{U}_k^{(j)} (\hat{U}_k^{(j)})^\top - \hat{U}_k^{(j-1)} (\hat{U}_k^{(j-1)})^\top\|_S \leq \epsilon,$$

- 6: Estimate and output:

$$\hat{U}_k^{\text{iFinal}} = \hat{U}_k^{(j)}, \quad k = 1, \dots, K,$$

$$\hat{P}_k^{\text{iFinal}} = \hat{U}_k^{\text{iFinal}} (\hat{U}_k^{\text{iFinal}})^\top, \quad k = 1, \dots, K,$$

$$\hat{\mathcal{F}}_t^{\text{iFinal}} = \mathcal{X}_t \times_{k=1}^K (\hat{U}_k^{\text{iFinal}})^\top, \quad t = 1, \dots, T,$$

$$\hat{\mathcal{E}}_t^{\text{iFinal}} = \mathcal{X}_t - \mathcal{X}_t \times_1 \hat{P}_1^{\text{iFinal}} \times_2 \dots \times_K \hat{P}_K^{\text{iFinal}}, \quad t = 1, \dots, T.$$

When we use the UTOPUP operator (2.5) for both UINIT and UITER in Algorithm 1, it will be called iTOPUP procedure. Similarly, iTIPUP uses UTIPUP operator (2.8) for both UINIT and UITER. Besides these two versions, we may also use UTIPUP for UINIT and UTOPUP for UITER, named as TIPUP-iTOPUP. Similarly, TOPUP-iTIPUP uses UTOPUP as UINIT and UTIPUP as UITER. These variants are sometimes useful, because TOPUP and TIPUP have different theoretical properties as the initializer or for iteration, as we will discuss in Section 3. Other estimators of the loading spaces based on the tensor time series can also be used in place of UINIT and UITER, such as the conventional high order SVD for tensor decomposition, which we refer to as Unfolding Procedure (UP), that simply performs SVD of the matricization along the appropriate mode of the $K + 1$ order tensor $(\mathcal{X}_1, \dots, \mathcal{X}_T)$ with time dimension as the additional $(K + 1)$ -th mode.

REMARK 2.1. While Algorithm 1 resembles an HOOI-type iteration of the orthogonal projection and singular matrix estimation methods, the proposed iTOPUP and iTIPUP are

significantly different from HOOI which iterates the operations of

orthogonal projection \rightarrow matrix unfolding \rightarrow SVD.

In both iTOPUP and iTIPUP, each iteration carries out the operations

(2.10) orthogonal projection \rightarrow autocovariance \rightarrow matrix unfolding \rightarrow SVD.

As the outer product is taken with TOPUP_k in (2.4), its orthogonal projection and autocovariance operations are exchangeable, so that we can write

$$\text{iTOPUP} = \text{HOOI}(\hat{\Sigma}_h, h = 1, \dots, h_0)$$

as long as the HOOI is modified by applying $U_\ell^{(j)}$ to both mode ℓ and mode $K + \ell, \ell \neq k$ in the projection operation and leaving alone the $(2K + 1)$ -th mode in the lags $1 : h_0$ throughout. However, for iTIPUP, the orthogonal projection and autocovariance operations in (2.10) are not exchangeable as the projections are sandwiched inside the autocovariance. Needless to say, the analysis of iTOPUP and iTIPUP is much more difficult than the conventional HOOI with iid assumption due to the involvement of the autocovariance operations in the time-axis in the iterations.

REMARK 2.2 (Rank determination). Here the estimators are constructed with given ranks r_1, \dots, r_K , though in theoretical analysis they are allowed to diverge. In practice, existing procedures for rank determination in the vector factor model, including the information criteria approach (Bai and Ng, 2002, 2007, Hallin and Liška, 2007) and ratio of eigenvalues approach (Lam and Yao, 2012, Ahn and Horenstein, 2013) can be extended to the tensor factor model by treating $d_1 \times \dots \times d_k$ tensors as d -dimensional vectors, $d = \prod_{k=1}^K d_k$.

3. Theoretical Properties. In this section we present some theoretical properties of the iterative procedures. We first present the additional notations needed for the discussion, then the error bounds for the iterative estimators under a minimum condition on the error process \mathcal{E}_t in the model. These error bounds are quite general and cover many different models. To help decipher the general results, we then present two concrete models (or still general sets of assumptions) of the signal process of the model, under which we will be able to obtain simpler and more explicit convergence rates.

3.1. Notations. We introduce some notations first. Let $\bar{\mathbb{E}}(\cdot) = \mathbb{E}(\cdot | \{\mathcal{F}_1, \dots, \mathcal{F}_T\})$. Define $d = \prod_{k=1}^K d_k$, $d_{-k} = d/d_k$, $r = \prod_{k=1}^K r_k$ and $r_{-k} = r/r_k$. Define order-4 tensors

$$(3.1) \quad \begin{aligned} \Theta_{k,h} &= \sum_{t=h+1}^T \frac{\text{mat}_k(\mathcal{M}_{t-h}) \otimes \text{mat}_k(\mathcal{M}_t)}{T-h} \in \mathbb{R}^{d_k \times d_{-k} \times d_k \times d_{-k}}, \\ \Phi_{k,h} &= \sum_{t=h+1}^T \frac{\text{mat}_k(\mathcal{F}_{t-h}) \otimes \text{mat}_k(\mathcal{F}_t)}{T-h} \in \mathbb{R}^{r_k \times r_{-k} \times r_k \times r_{-k}}, \\ \Phi_{k,h}^{(\text{cano})} &= \sum_{t=h+1}^T \frac{\text{mat}_k(\mathcal{M}_{t-h} \times_{k=1}^K U_k^\top) \otimes \text{mat}_k(\mathcal{M}_t \times_{k=1}^K U_k^\top)}{T-h} \in \mathbb{R}^{r_k \times r_{-k} \times r_k \times r_{-k}}, \end{aligned}$$

with U_k from the SVD $A_k = U_k \Lambda_k V_k^\top$. We view $\Phi_{k,h}^{(\text{cano})}$ as the canonical version of the auto-covariance of the factor process. The noiseless version of the matrix TOPUP_k in (2.4) is

$$(3.2) \quad \text{mat}_1(\Theta_{k,1:h_0}) = \bar{\mathbb{E}}[\text{TOPUP}_k] \in \mathbb{R}^{d_k \times (dd_{-k}h_0)},$$

with $\Theta_{k,1:h_0} = (\Theta_{k,h}, h = 1, \dots, h_0)$. The canonical factor version of (3.2) is $\text{mat}_1(\Phi_{k,1:h_0}^{(\text{cano})}) \in \mathbb{R}^{r_k \times (r_{r-k} h_0)}$ with $\Phi_{k,1:h_0}^{(\text{cano})} = (\Phi_{k,h}^{(\text{cano})}, h = 1, \dots, h_0) \in \mathbb{R}^{r_k \times r_{r-k} \times r_k \times r_{r-k} \times h_0}$. Similarly define

$$(3.3) \quad \begin{aligned} \Theta_{k,h}^* &= \sum_{t=h+1}^T \frac{\text{mat}_k(\mathcal{M}_{t-h}) \text{mat}_k^\top(\mathcal{M}_t)}{T-h} \in \mathbb{R}^{d_k \times d_k}, \\ \Phi_{k,h}^* &= \sum_{t=h+1}^T \frac{\text{mat}_k(\mathcal{F}_{t-h}) \text{mat}_k^\top(\mathcal{F}_t)}{T-h} \in \mathbb{R}^{r_k \times r_k}, \\ \Phi_{k,h}^{*(\text{cano})} &= U_k^\top \Theta_{k,h}^* U_k \\ &= \sum_{t=h+1}^T \frac{\text{mat}_k(\mathcal{M}_{t-h} \times_{k=1}^K U_k^\top) \text{mat}_k^\top(\mathcal{M}_t \times_{k=1}^K U_k^\top)}{T-h} \in \mathbb{R}^{r_k \times r_k}. \end{aligned}$$

The noiseless version of (2.7) is

$$(3.4) \quad \Theta_{k,1:h_0}^* = (\Theta_{k,h}^*, h = 1, \dots, h_0) = \overline{\mathbb{E}}[\text{TIPUP}_k] \in \mathbb{R}^{d_k \times (d_k h_0)}$$

and its canonical factor version is $\Phi_{k,1:h_0}^{*(\text{cano})} = (\Phi_{k,h}^{*(\text{cano})}, h = 1, \dots, h_0) \in \mathbb{R}^{r_k \times (r_k h_0)}$. Let $\tau_{k,m}$ be the m -th singular value of the noiseless version of the TOPUP $_k$ matrix,

$$\tau_{k,m} = \sigma_m(\overline{\mathbb{E}}[\text{TOPUP}_k]) = \sigma_m(\text{mat}_1(\Theta_{k,1:h_0}^*)) = \sigma_m(\text{mat}_1(\Phi_{k,1:h_0}^{*(\text{cano})})).$$

The signal strength for iTOPUP can be characterized as

$$(3.5) \quad \lambda_k = \sqrt{h_0^{-1/2} \tau_{k,r_k}}.$$

Similarly, let

$$\tau_{k,m}^* = \sigma_m(\overline{\mathbb{E}}(\text{TIPUP}_k)) = \sigma_m(\Theta_{k,1:h_0}^*) = \sigma_m(\Phi_{k,1:h_0}^{*(\text{cano})}).$$

The signal strength for iTIPUP can be characterized as

$$(3.6) \quad \lambda_k^* = \sqrt{h_0^{-1/2} \tau_{k,r_k}^*}.$$

We note that by (3.3) and the Cauchy-Schwarz inequality,

$$\lambda_k^{*2} \leq h_0^{-1/2} \|\Theta_{k,1:h_0}^*\|_S \leq \max_{h \leq h_0} \|\Theta_{k,h}^*\|_S \leq \|\Theta_{k,0}^*\|_S / (1 - h_0/T).$$

3.2. General error bounds. Our general error bounds for the proposed iTOPUP and iTIPUP are established under the following assumption for the error process.

ASSUMPTION 1. The error process \mathcal{E}_t are independent Gaussian tensors conditionally on the factor process $\{\mathcal{F}_t, t \in \mathbb{Z}\}$. In addition, there exists some constant $\sigma > 0$, such that

$$\overline{\mathbb{E}}(u^\top \text{vec}(\mathcal{E}_t))^2 \leq \sigma^2 \|u\|_2^2, \quad u \in \mathbb{R}^d.$$

Assumption 1 is used by [Chen, Yang and Zhang \(2019\)](#) for the theoretical investigation of the non-iterative TIPUP and TOPUP, and is similar to those on the noise imposed in [Lam, Yao and Bathia \(2011\)](#), [Lam and Yao \(2012\)](#). The normality assumption, which ensures fast convergence rates in our analysis, is imposed for technical convenience. It accommodates general patterns of dependence among individual time series fibers, but also allows a presentation of the main results with manageable analytical complexity. In fact, direct extension is

visible in our analysis under the sub-Gaussian and even more general tail probability conditions. Under Assumption 1 the magnitude of the noise can be measured by the dimension d_k before the projection and by the rank r_k after the projection. The main theorems (Theorems 3.1, 3.2 and 3.3) in this section are based on this assumption on the noise alone, and cover all thereafter discussed settings of the signal \mathcal{M}_t .

Let us first study the behavior of iTOPUP procedure. By Chen, Yang and Zhang (2019), the risk $\mathbb{E}[\|\widehat{U}_k^{(0)}\widehat{U}_k^{(0)\top} - U_k U_k^\top\|_S]$ of the TOPUP estimator for U_k , the initialization of iTOPUP, is no larger than a constant times

$$(3.7) \quad R_k^{(0)} = \lambda_k^{-2} \sigma T^{-1/2} \left\{ \sqrt{d_k d_{-k} r_{-k}} \|\Theta_{k,0}^*\|_S^{1/2} + (\sqrt{d_k} + \sqrt{d_{-k} r_{-k}}) \|\Theta_{k,0}\|_{\text{op}}^{1/2} + \sigma \sqrt{d_k d_{-k}} + \sigma d_k \sqrt{d_{-k}} T^{-1/2} \right\},$$

where $d_{-k} = \prod_{j \neq k} d_j$ and $r_{-k} = \prod_{j \neq k} r_j$. The aim of iTOPUP is to achieve dimension reduction by projecting the data in other modes of the tensor time series from \mathbb{R}^{d_j} to \mathbb{R}^{r_j} , $j \neq k$. Ideally (e.g. when the true projection matrices U_j are used), this would reduce the above rate in (3.7) to

$$(3.8) \quad R_k^{(\text{ideal})} = \lambda_k^{-2} \sigma T^{-1/2} \left\{ \sqrt{d_k r_{-k}} \|\Theta_{k,0}^*\|_S^{1/2} + (\sqrt{d_k} + \sqrt{r_{-k} r}) \|\Theta_{k,0}\|_{\text{op}}^{1/2} + \sigma \sqrt{d_k r_{-k}} + \sigma d_k \sqrt{r_{-k}} T^{-1/2} \right\},$$

by replacing all d_j with r_j , $j \neq k$. However, because the iteration uses the estimated U_j , $j \neq k$, of total dimension $d_{-k}^* = \sum_{j \neq k} d_j r_j$, our analysis also involves the following additional error term,

$$(3.9) \quad R_k^{(\text{add})} = \lambda_k^{-2} \sigma^2 T^{-1} \left(d_{-k}^* + \sqrt{d_{-k}^* d_k r_{-k}} \right).$$

The following theorem provides conditions under which the ideal rate is indeed achieved.

THEOREM 3.1. *Suppose Assumption 1 holds. Let $h_0 \leq T/4$ and P_k , $\Theta_{k,0}$, $\Theta_{k,0}^*$ and λ_k be as in (2.2), (3.1), (3.3) and (3.5) respectively. Let $R^{(0)} = \max_{1 \leq k \leq K} R_k^{(0)}$ with the $R_k^{(0)}$ in (3.7), $R^{(\text{ideal})} = \max_{1 \leq k \leq K} R_k^{(\text{ideal})}$ with the $R_k^{(\text{ideal})}$ in (3.8), and $R^{(\text{add})} = \max_{1 \leq k \leq K} R_k^{(\text{add})}$ with the $R_k^{(\text{add})}$ in (3.9). Let $\widehat{P}_k^{(m)} = \widehat{U}_k^{(m)} \widehat{U}_k^{(m)\top}$ with the m -step estimator $\widehat{U}_k^{(m)}$ in the iTOPUP algorithm. Then, the following statements hold for a certain numerical constant $C_1^{(\text{TOPUP})}$ and a constant $C_{1,K}^{(\text{iter})}$ depending on K only: When*

$$(3.10) \quad C_1^{(\text{TOPUP})} R^{(0)} \leq (1 - \rho)/4 \text{ and } C_{1,K}^{(\text{iter})} (R^{(\text{ideal})} + R^{(\text{add})}) \leq \rho$$

with a constant $0 < \rho < 1$, it holds simultaneously for all $1 \leq k \leq K$ and $m \geq 0$ that

$$(3.11) \quad \|\widehat{P}_k^{(m)} - P_k\|_S \leq 2C_1^{(\text{TOPUP})} \left((1 - \rho^m)(1 - \rho)^{-1} R^{(\text{ideal})} + (\rho^m/2) R^{(0)} \right)$$

in an event with probability at least $1 - \sum_{k=1}^K e^{-d_k}$. In particular, after at most $J = \lceil \log(\max_k d_{-k}/r_{-k}) / \log(1/\rho) \rceil$ iterations,

$$(3.12) \quad \mathbb{E} \left[\max_{1 \leq k \leq K} \|\widehat{P}_k^{(J)} - P_k\|_S \right] \leq \frac{3C_1^{(\text{TOPUP})}}{1 - \rho} R^{(\text{ideal})} + \sum_{k=1}^K e^{-d_k}.$$

REMARK 3.1. The essence of our analysis of iTOPUP is that under (3.10), each iteration is a contraction of the error in the estimation of $\times_{j \neq k} U_j$ in a small neighborhood of it. The upper bound (3.11) for the error of the m -step estimator is comprised of two terms respectively corresponding to the cumulative iteration error and the contracted error of the initial estimator. Of course, after sufficiently large number of iterations, the first term would dominate the second as in (3.12).

REMARK 3.2. The constant $C_1^{(\text{TOPUP})}$ is taken in (3.10) to guarantee sufficient accuracy of the initialization of iTOPUP in the following sense:

$$(3.13) \quad \max_{k \leq K} \mathbb{E} \|\widehat{U}_k^{(0)} (\widehat{U}_k^{(0)})^\top - P_k\|_{\text{S}} \leq C_1^{(\text{TOPUP})} R^{(0)}$$

with at least probability $1 - 5^{-1} \sum_{k=1}^K e^{-d_k}$. The consistency of the non-iterative TOPUP estimator requires $R^{(0)} \rightarrow 0$ (Chen, Yang and Zhang, 2019). However, here we do not require the TOPUP estimator as the initial value to be consistent. For (3.11) to hold, the TOPUP estimator is only required to be sufficiently close to the ground truth as in (3.13).

REMARK 3.3. It is relatively easy to verify that the first part of (3.10) implies the second part under many circumstances, including when d_k are of the same order, r_k are of the same order, and $r_k \lesssim d_k^{1-1/K}$ ($K \geq 2$). In Zhang and Xia (2018), condition $\max_k r_k \lesssim \min_k d_k^{1/2}$ is imposed to control the complexity of the estimated U_j in HOOI although their error bound is sharp and their model is very different. In Corollaries 3.1 and 3.3 below, we prove that the second part of (3.10) follows from the first part respectively in a general fixed rank model and a general diverging rank model. In fact $R_k^{(\text{ideal})} + R_k^{(\text{add})} \ll R^{(0)}$ typically so that the second part of (3.10) provides a non-asymptotic lower bound for the ρ in (3.11), allowing $\rho = \rho_{T, d_k, d_{-k}^*, r_k, r_{-k}, \lambda_k} \rightarrow 0$. In Corollary 3.1 below, $\rho = C_{1, K}^{(\text{iter})} (R^{(\text{ideal})} + R^{(\text{add})})$ is taken in (3.10) to give (3.12) in one iteration when $R_k^{(\text{ideal})}$ dominates $R_k^{(\text{add})}$.

REMARK 3.4. When the loading matrices A_k and the TOPUP version of the matrix unfolding of the auto-covariance of \mathcal{F}_t all have bounded condition numbers and average squared entries of magnitude 1, λ_k^2 , $\|\Theta_{k,0}^*\|_{\text{S}}$ and $\|\Theta_{k,0}\|_{\text{op}}$ are all of the order $d \times \text{poly}(r_1, \dots, r_K)$. In this case, Theorem 3.1 just requires $T \geq \text{poly}(r_1, \dots, r_K)$ for the initialization to achieve through iteration the fast convergence rate $T^{-1/2} d_{-k}^{-1/2} \text{poly}(r_1, \dots, r_K)$. See Corollary 3.3 for details. This is in sharp contrast to the results of traditional factor analysis which requires $T \rightarrow \infty$ to consistently estimate the loading spaces. The main reason is that the other tensor modes provide additional information and in certain sense serve as additional samples. Roughly speaking, we have totally $dT = d_k d_{-k} T$ observations in the tensor time series to estimate the $d_k r_k$ parameters in the projection to the column space of the loading matrix A_k , where $r_k \ll d_{-k} T$ in the above ‘‘regular’’ case.

Now, let us consider the statistical performance of iTIPUP procedure. Again, by Chen, Yang and Zhang (2019) the TIPUP risk in the estimation of P_k is bounded by

$$(3.14) \quad \mathbb{E} [\|\widehat{P}_k^{(\text{TIPUP})} - P_k\|_{\text{S}}] \lesssim R_k^{*(0)} = (\lambda_k^*)^{-2} \sigma T^{-1/2} \sqrt{d_k} \left(\|\Theta_{k,0}^*\|_{\text{S}}^{1/2} + \sigma \sqrt{d_{-k}} \right)$$

with $d_{-k} = \prod_{j \neq k} d_j$, and the aim of iTIPUP is to achieve the ideal rate

$$(3.15) \quad R_k^{*(\text{ideal})} = (\lambda_k^*)^{-2} \sigma T^{-1/2} \sqrt{d_k} \left(\|\Theta_{k,0}^*\|_{\text{S}}^{1/2} + \sigma \sqrt{r_{-k}} \right)$$

through dimension reduction, where $r_{-k} = \prod_{j \neq k} r_j$. As in the case of iTOPUP, our error bound for iTIPUP involves the additional error term

$$(3.16) \quad R_k^{*(\text{add})} = \sqrt{d_{-k}^*/d_k} R_k^{*(\text{ideal})}.$$

The following theorem, which allows the ranks r_k to grow to infinity as well as d_k when $T \rightarrow \infty$, provides sufficient conditions to guarantee the ideal convergence rate for iTIPUP.

THEOREM 3.2. *Suppose Assumption 1 holds. Let P_k , $\Theta_{k,0}^*$ and λ_k^* be as in (2.2), (3.3) and (3.6) respectively. Let $h_0 \leq T/4$, and*

$$R^{*(0)} = \max_{1 \leq k \leq K} R_k^{*(0)}; \quad R^{*(\text{ideal})} = \max_{1 \leq k \leq K} R_k^{*(\text{ideal})}, \quad R^{*(\text{add})} = \max_{1 \leq k \leq K} R_k^{*(\text{add})}.$$

with $R_k^{*(0)}$ in (3.14), $R_k^{*(\text{ideal})}$ in (3.15) and $R_k^{*(\text{add})}$ in (3.16). Let $\hat{P}_k^{(m)} = \hat{U}_k^{(m)} \hat{U}_k^{(m)\top}$ with the m -step estimator $\hat{U}_k^{(m)}$ in iTIPUP algorithm. Then, the following statements hold for a certain numerical constant $C_1^{(\text{TIPUP})}$ and a constant $C_{1,K}^{(\text{iter})}$ depending on K only: When

$$(3.17) \quad C_1^{(\text{TIPUP})} R^{*(0)} \leq (1 - \rho) \frac{\min_{1 \leq k \leq K} \lambda_k^{*2}}{8 \|\Theta_{k,0}^*\|_S} \quad \text{and} \quad C_{1,K}^{(\text{iter})} (R^{*(\text{ideal})} + R^{*(\text{add})}) \leq \rho$$

with a constant $0 < \rho < 1$, it holds simultaneously for all $1 \leq k \leq K$ and $m \geq 0$ that

$$(3.18) \quad \|\hat{P}_k^{(m)} - P_k\|_S \leq 2C_1^{(\text{TIPUP})} \left((1 - \rho^m)(1 - \rho)^{-1} R^{*(\text{ideal})} + (\rho^m/2) R^{*(0)} \right)$$

in an event with probability at least $1 - \sum_{k=1}^K e^{-d_k}$. In particular, after at most $J = \lceil \log(\max_k d_{-k}/r_{-k}) / \log(1/\rho) \rceil$ iterations,

$$(3.19) \quad \mathbb{E} \left[\max_{1 \leq k \leq K} \|\hat{P}_k^{(J)} - P_k\|_S \right] \leq \frac{3C_1^{(\text{TIPUP})}}{1 - \rho} R^{*(\text{ideal})} + \sum_{k=1}^K e^{-d_k}.$$

We briefly discuss the conditions and conclusions of Theorem 3.2 as the details are parallel to the remarks below Theorem 3.1. By (3.3), (3.6) and the Cauchy-Schwarz inequality, $(1 - h_0/T) \lambda_k^{*2} \leq \|\Theta_{k,0}^*\|_S$, so that the first condition in (3.17) guarantees a sufficiently small $R^{*(0)}$, which implies a sufficiently small error in the initialization of iTIPUP by (3.14). The second condition in (3.17) again has two terms respectively reflecting the ideal rate after dimension reduction by the true $U_{-k} = \odot_{j \neq k} U_j$ in the estimation of U_k and the extra cost of estimating U_{-k} . The upper bound (3.18) for the error of the m -step estimator is also comprised of two terms representing the cumulative iteration error and contracted initialization error. In Corollary 3.2 below with fixed r_k , the smallest $\rho = C_{1,K}^{(\text{iter})} (R^{*(\text{ideal})} + R^{*(\text{add})})$ is taken in (3.17) to achieve (3.19) in one iteration when $R_k^{*(\text{ideal})}$ dominates $R_k^{*(\text{add})}$. Moreover, Theorem 3.2 allows diverging ranks r_k and convergence rate $T^{-1/2} d_{-k}^{-1/2} \text{poly}(r_1, \dots, r_K)$ under proper conditions as discussed in Remark 3.4.

As discussed in Section 2.3, we can mix the TOPUP and TIPUP operations for the initiation and iterative operations in Algorithm 1. For example, the proof of Theorems 3.1 yields the following error bound for the mixed TIPUP-iTOPUP algorithm.

THEOREM 3.3. *Assumption 1 holds. Let $R^{(0)}$, $R^{(\text{ideal})}$ and $R^{(\text{add})}$ be as in Theorem 3.1 and $R^{*(0)}$ be as in Theorem 3.2. Let $\hat{P}_k^{(m)} = \hat{U}_k^{(m)} \hat{U}_k^{(m)\top}$ with $\hat{U}_k^{(m)}$ being the m -step estimator in the TIPUP-iTOPUP algorithm. Then, the following statement holds for a certain numerical constant $C_1^{(\text{TOPUP})}$ and a constant $C_{1,K}^{(\text{iter})}$ depending on K only: When*

$$(3.20) \quad C_1^{(\text{TOPUP})} R^{*(0)} \leq (1 - \rho)/4 \quad \text{and} \quad C_{1,K}^{(\text{iter})} (R^{(\text{ideal})} + R^{(\text{add})}) \leq \rho$$

with a constant $0 < \rho < 1$, it holds in an event with probability at least $1 - \sum_{k=1}^K e^{-d_k}$ that simultaneously for all $1 \leq k \leq K$ and $m \geq 0$

$$\|\widehat{P}_k^{(m)} - P_k\|_S \leq 2C_1^{(\text{TOPUP})} \left((1 - \rho^m)(1 - \rho)^{-1} R^{(\text{ideal})} + (\rho^m/2) R^{*(0)} \right).$$

We omit the statement of an analogous error bound for the TOPUP-iTIPUP algorithm.

3.3. Fixed rank factor process. In this section we provide the convergence rate when the dimensions of the factors \mathcal{F}_t , or equivalently the ranks of the signal process $\mathcal{M}_t, r_1, \dots, r_K$, are fixed, and the auto-cross-outer-product of the factor process is ergodic. Formally, we impose the following additional assumption.

ASSUMPTION 2. The ranks r_1, \dots, r_K are fixed. The factor process \mathcal{F}_t is weakly stationary and its auto-cross-outer-product process is ergodic in the sense of

$$\frac{1}{T-h} \sum_{t=h+1}^T \mathcal{F}_{t-h} \otimes \mathcal{F}_t \longrightarrow \mathbb{E}(\mathcal{F}_{t-h} \otimes \mathcal{F}_t) \quad \text{in probability,}$$

where the elements of $\mathbb{E}(\mathcal{F}_{t-h} \otimes \mathcal{F}_t)$ are all finite. In addition, the condition numbers of $A_k^\top A_k$ ($k = 1, \dots, K$) are bounded. Furthermore, assume that h_0 is fixed, and

- (i) (TOPUP related): $\mathbb{E}[\text{mat}_1(\Phi_{k,1:h_0})]$ is of rank r_k for $1 \leq k \leq K$.
- (ii) (TIPUP related): $\mathbb{E}[\text{mat}_1(\Phi_{k,1:h_0}^{*(\text{cano})})]$ is of rank r_k for $1 \leq k \leq K$.

Under Assumption 2, the factor process has a fixed expected auto-cross-moment tensor with fixed dimensions. The assumption that the condition numbers of $A_k^\top A_k$ ($k = 1, \dots, K$) are bounded corresponds to the pervasive condition (e.g., [Stock and Watson \(2002\)](#), [Bai \(2003\)](#)). It ensures that all the singular values of A_k are of the same order. Such conditions are commonly imposed in factor analysis.

As our methods are based on auto-cross-moment at nonzero lags, we do not need to assume any specific model for the latent process \mathcal{F}_t , except some rank conditions in Assumption 2(i) and (ii). Since the columns of $\text{mat}_1(\Phi_{k,1:h_0}^{*(\text{cano})})$ are linear combinations of those of $\text{mat}_1(\Phi_{k,1:h_0}^{(\text{cano})})$ and $\mathbb{E}[\text{mat}_1(\Phi_{k,1:h_0})]$ and $\mathbb{E}[\text{mat}_1(\Phi_{k,1:h_0}^{(\text{cano})})]$ have the same rank, Assumption 2(ii) implies Assumption 2(i).

In order to provide a more concrete understanding of Assumption 2(i) and (ii), consider the case of $k = 1$ and $K = 2$. We write the factor process $\mathcal{F}_t = (f_{i,j,t})_{r_1 \times r_2}$, and the stationary auto-cross-moments $\phi_{i_1, j_1, i_2, j_2, h} = \mathbb{E}(f_{i_1, j_1, t-h} f_{i_2, j_2, t})$. Hence $\mathbb{E}[\text{mat}_1(\Phi_{k,1:h_0})]$ is a $r_k \times (r_{-k} r_k r_{-k} h_0)$ matrix, with columns being $\phi_{\cdot, j_1, i_2, j_2, h}$. Since $\mathbb{E}[\text{mat}_1(\Phi_{k,1:h_0})] \mathbb{E}[\text{mat}_1(\Phi_{k,1:h_0})]^\top$ is a sum of many semi-positive definite $r_k \times r_k$ matrices, if any one of these matrices is full rank, then $\mathbb{E}[\text{mat}_1(\Phi_{k,1:h_0})]$ is of rank r_k . Hence Assumption 2(i) is relatively easy to fulfill. On the other hand, Assumption 2(ii) is quite different. First, the condition is imposed on the canonical form of the model as the inner product in TIPUP related procedures behaves differently. Let $\mathcal{F}_t^{(\text{cano})} = U_1^\top \mathcal{M}_t U_2 = (f_{i,j,t}^{(\text{cano})})_{r_1 \times r_2}$, and $\phi_{i_1, j_1, i_2, j_2, h}^{(\text{cano})} = \mathbb{E}(f_{i_1, j_1, t-h}^{(\text{cano})} f_{i_2, j_2, t}^{(\text{cano})})$. Then $\|\Phi_{1,1:h_0}^{*(\text{cano})}\|_{\text{HS}}^2 = \sum_{h=1}^{h_0} \sum_{i_1, i_2} \left(\sum_{j=1}^{r_2} \phi_{i_1, j, i_2, j, h}^{(\text{cano})} \right)^2$. As $\phi_{i_1, j, i_2, j, h}^{(\text{cano})}$ may be positive or negative for different i_1, i_2, j, h , the summation $\sum_{j=1}^{r_2} \phi_{i_1, j, i_2, j, h}^{(\text{cano})}$ is subject to potential signal cancellation for $h > 0$. Assumption 2(ii) ensures that there is no complete signal cancellation that makes the rank of $\mathbb{E}[\text{mat}_1(\Phi_{k,1:h_0}^{*(\text{cano})})]$ less than r_k . While the signal cancellation rarely causes the rank deficiency, the resulting loss of efficiency may still have an impact on the finite sample performance as our simulation results demonstrate. Of course complete signal cancellation is less likely with larger h_0 .

The following corollary is a simplified version of Theorem 3.1 under Assumption 2(i).

COROLLARY 3.1. *Suppose Assumptions 1 and 2(i) hold. Let $\lambda = \prod_{k=1}^K \|A_k\|_S$ and $d_{\min} = \min\{d_1, \dots, d_K\}$. Let $h_0 \leq T/4$ and σ fixed. Then, there exist numerical constants $C_{0,K,r}$ and $C_{1,K,r}$ depending on K and r only such that when*

$$(3.21) \quad \lambda^2 \geq C_{0,K,r} \sigma^2 \left(\frac{d}{T} + \frac{d}{\sqrt{T d_{\min}}} \right),$$

the 1-step iTOPUP estimator satisfies

$$(3.22) \quad \mathbb{E} \|\hat{P}_k^{(1)} - P_k\|_S \leq C_{1,K,r} \left(\frac{\sigma \sqrt{d_k}}{\lambda \sqrt{T}} + \frac{\sigma^2 \sqrt{d_k}}{\lambda^2 \sqrt{T}} \right) + \sum_{k=1}^K e^{-d_k}.$$

Corollary 3.1 asserts that, in order to recover the factor loading space for A_k , the signal to noise ratio needs to satisfy $\lambda/\sigma \geq C_0(d^{1/2}T^{-1/2} + d^{1/2}d_{\min}^{-1/4}T^{-1/4})$ as in (3.21), and the ideal rate (3.22) can be achieved in one iteration. The ideal rate is much sharper than the convergence rate of the non-iterative TOPUP in Chen, Yang and Zhang (2019).

The following corollary is a simplified version of Theorem 3.2 under Assumption 2(ii), which excludes severe signal cancellation in iTIPUP.

COROLLARY 3.2. *Suppose Assumptions 1 and 2(ii) hold. Let $\lambda = \prod_{k=1}^K \|A_k\|_S$ and $d_{\max} = \max\{d_1, \dots, d_K\}$. Let $h_0 \leq T/4$ and σ fixed. Then, there exist constants $C_{0,K,r}$ and $C_{1,K,r}$ depending on K and r only such that when*

$$(3.23) \quad \lambda^2 \geq C_{0,K,r} \sigma^2 \left(\frac{d_{\max}}{T} + \sqrt{\frac{d}{T}} \right),$$

both the 1-step iTIPUP estimator and the 1-step TIPUP-iTOPUP estimator satisfy

$$(3.24) \quad \mathbb{E} \|\hat{P}_k^{(1)} - P_k\|_S \leq C_{1,K,r} \left(\frac{\sigma \sqrt{d_k}}{\lambda \sqrt{T}} + \frac{\sigma^2 \sqrt{d_k}}{\lambda^2 \sqrt{T}} \right) + \sum_{k=1}^K e^{-d_k}.$$

Compared with the results in Corollary 3.1 for iTOPUP, the achieved ideal rate (3.24) is the same. However, the signal-to-noise ratio requirement (3.23) is weaker but Assumption 2(ii) is stronger in Corollary 3.2 for iTIPUP. Again, the ideal rate is much sharper than the convergence rate of the non-iterative TIPUP in Chen, Yang and Zhang (2019).

3.4. Diverging ranks. The main theorems in Subsection 3.2 allow for the case where the dimensions of the core factor, r_1, \dots, r_K , diverge as the dimensions of the observed tensor d_1, \dots, d_K grow to infinity. The following assumption provides a concrete set of conditions that can be used to provide some insights of the properties of iTOPUP and iTIPUP in such scenarios.

ASSUMPTION 3. For a certain $\delta_0 \in [0, 1]$, $\|\Theta_{k,0}\|_{\text{op}} \asymp \sigma^2 d^{1-\delta_0}/r$ and $\|\Theta_{k,0}^*\|_S \asymp \sigma^2 d^{1-\delta_0}/r_k$ with probability approaching one. For the singular values, two scenarios are considered.

- (i) (TOPUP related): There exist some constants $\delta_1 \in [\delta_0, 1]$ and $c_1 > 0$ such that with probability approaching one (as $T \rightarrow \infty$) $\lambda_k^2 \geq c_1 \sigma^2 d^{1-\delta_1} / \sqrt{r r_k}$, for all $k = 1, \dots, K$.
- (ii) (TIPUP related): There exist some constants $\delta_1 \in [\delta_0, 1]$, $c_2 > 0$ and $\delta_2 \geq 0$ such that with probability approaching one (as $T \rightarrow \infty$), $\lambda_k^{*2} \geq c_2 \sigma^2 d^{1-\delta_1} r_k^{-1} r_{-k}^{-\delta_2}$ for all $k = 1, \dots, K$.

Assumption 3 is similar to the signal strength condition of Lam and Yao (2012), and the pervasive condition on the factor loadings (e.g., Stock and Watson (2002) and Bai (2003)). It is more general than Assumption 2 in the sense that it allows r_1, \dots, r_K to diverge and the latent process \mathcal{F}_t does not have to be weakly stationary.

We take δ_0, δ_1 as measures of the strength of the signal process \mathcal{M}_t . They roughly indicate how much information is contained in the signals compared with the amount of noise, with respect to the dimensions and ranks, d, r and r_k . In this sense, they reflect the signal to noise ratio. When $\delta_0 = \delta_1 = 0$, the factors are called strong factors; otherwise, the factors are called weak factors.

REMARK 3.5 (Signal Strength and the index δ_0). We note that $\text{trace}(\Theta_{k,0}) = \text{trace}(\Theta_{k,0}^*) = \sum_{t=1}^T \|\text{vec}(\mathcal{M}_t)\|_2^2 / T$, and that $\text{rank}(\Theta_{k,0}) = r$ and $\text{rank}(\Theta_{k,0}^*) = r_k$ when the data is in general position, where $\Theta_{k,0}$ is treated as a $d \times d$ matrix. Thus, if $\sum_{t=1}^T \|\text{vec}(\mathcal{M}_t)\|_2^2 / (\sigma^2 d T) \asymp d^{-\delta_0}$ is the signal-to-noise ratio, then the condition $\|\Theta_{k,0}\|_{\text{op}} \asymp \sigma^2 d^{1-\delta_0} / r$ holds when r is the order of the effective rank of $\Theta_{k,0}$ and the condition $\|\Theta_{k,0}^*\|_{\text{S}} \asymp \sigma^2 d^{1-\delta_0} / r_k$ holds when r_k is the order of the effective rank of $\Theta_{k,0}^*$. Because the signal \mathcal{M}_t has d elements at each t , the assumption $\sum_{t=1}^T \|\text{vec}(\mathcal{M}_t)\|_2^2 / (\sigma^2 d T) \asymp d^{-\delta_0}$ says that the squared ratio of the elements and the noise level is $d^{-\delta_0}$ averaged over time and space. Thus, the factor is called strong when $\delta_0 = 0$. In view of (1.1) and (1.2), $\mathcal{M}_t = \mathcal{F}_t \times_{k=1}^K A_k$, so that we may have weaker factor with $\delta_0 > 0$ when the loading matrices A_k are sparse or have some relatively small singular components. We note that by Cauchy-Schwarz, the signal-to-noise ratio conditions also imply $(1 - h/T)^2 \|\Theta_{k,h}\|_{\text{HS}}^2 \leq \|\Theta_{k,0}\|_{\text{HS}}^2 \lesssim r(\sigma^2 d^{1-\delta_0} / r)^2$ and $(1 - h/T)^2 \|\Theta_{k,h}^*\|_{\text{HS}}^2 \leq \|\Theta_{k,0}^*\|_{\text{HS}}^2 \lesssim r_k(\sigma^2 d^{1-\delta_0} / r_k)^2$ respectively.

REMARK 3.6 (Assumption 3(i) and the role of δ_1). In fact, for TOPUP, Assumption 3(i) holds when (a) $\|\mathbb{E}[\text{TOPUP}_k]\|_{\text{HS}}^2 = \sum_{h=1}^{h_0} \|\Theta_{k,h}\|_{\text{HS}}^2 \asymp h_0 \sigma^4 d^{2(1-\delta_1)} / r$ and (b) all the nonzero singular values of $\mathbb{E}[\text{TOPUP}_k]$ are of the same order. Because $\|\Theta_{k,h}\|_{\text{HS}}^2 \lesssim \sigma^4 d^{2(1-\delta_0)} / r$ by the condition on the signal-to-noise ratio, we must have $\delta_1 \geq \delta_0$, and $d^{\delta_0 - \delta_1}$ can be viewed as the order of average auto-correlation over lags $h = 1, \dots, h_0$. For $k = 1$ and $K = 2$, the factor process in the canonical form is $\mathcal{F}_t^{(\text{cano})} = U_1^\top \mathcal{M}_t U_2 = (f_{i,j,t}^{(\text{cano})})_{r_1 \times r_2}$, and $\phi_{i_1, j_1, i_2, j_2, h}^{(\text{cano})} = \sum_{t=h+1}^T f_{i_1, j_1, t-h}^{(\text{cano})} f_{i_2, j_2, t}^{(\text{cano})} / (T - h)$ is the time average cross product between the factor fibers $f_{i_1, j_1, 1:T}^{(\text{cano})}$ and $f_{i_2, j_2, 1:T}^{(\text{cano})}$. Thus, the first condition (a) means $\sum_{h=1}^{h_0} \|\Theta_{1,h}\|_{\text{HS}}^2 = \sum_{h=1}^{h_0} \|\Phi_{1,h}^{(\text{cano})}\|_{\text{HS}}^2 = \sum_{i_1, j_1, i_2, j_2, h} (\phi_{i_1, j_1, i_2, j_2, h}^{(\text{cano})})^2 \asymp h_0 \sigma^4 d^{2(1-\delta_1)} / r$.

REMARK 3.7 (Assumption 3(ii), the role of δ_2 and signal cancellation). The points parallel to those in Remark 3.6 are applicable to TIPUP, but with one caveat: Beyond the average auto-correlation, an additional discount $r_{-k}^{-\delta_2} \leq 1$ is needed to take into account the impact of possible signal cancellation with TIPUP and its iteration. For $k = 1$ and $K = 2$, $\|\Theta_{1,h}^*\|_{\text{HS}}^2 = \|\Phi_{1,h}^{(\text{cano})}\|_{\text{HS}}^2 = \sum_{i_1, i_2} \left(\sum_{j=1}^{r_2} \phi_{i_1, j, i_2, j, h}^{(\text{cano})} \right)^2$, and the summation inside the square is subject to signal cancellation for $h > 0$ since the auto-cross-moment $\phi_{i_1, j, i_2, j, h}^{(\text{cano})}$ can have different signs. The additional parameter δ_2 measures the severity of signal cancellation in the TIPUP related procedures. For example, when the majority of $\phi_{i_1, j, i_2, j, h}^{(\text{cano})}$ are of the same sign for most of (i_1, i_2, h) , it would be reasonable to assume $\delta_2 = 0$. When $\phi_{i_1, j, i_2, j, h}^{(\text{cano})}$ behave like independent mean zero variables, δ_2 would be close to 0.5. And $\delta_2 = \infty$ when all the signals cancel out by the summation $\phi_{i_1, j, i_2, j, h}^{(\text{cano})}$ over j . In the case of fixed r_k , the convergence rate depends on whether $\delta_2 = \infty$ (severe signal cancellation) or not.

REMARK 3.8 (The role of h_0). The selection of h_0 is a relative minor problem in practice though very complex to analyze. Theoretically it suffices to use an h_0 with λ_k of the right order, so that choosing a somewhat large h_0 would not harm the convergence rate for the proposed methods. In practice a small h_0 (less than 3) is often sufficient. The impact of the choice of h_0 on the signal and noise depends on the autocorrelation of the factor process, as well as the loading matrices. For example, if the factor process is of very short memory (e.g. an MA(1) process), including any lag $h > 1$ only introduces noise to TOPUP $_k$ in (2.4) and TIPUP $_k$ in (2.7) without enhancing the signal. On the other hand, including an extra lag is the most simple and effective way to prevent signal cancellation with iTIPUP, as discussed in the previous remark. Increasing h_0 includes more non-negative terms in the signal strength $\sum_{i_1, i_2, h} \left(\sum_{j=1}^{r_2} \phi_{i_1, j, i_2, j, h}^{(\text{cano})} \right)^2$, hence potentially reducing the chance of severe signal cancellation. The simulation results presented in the supplementary material provide some empirical behavior of choosing different h_0 . While the choice of h_0 will affect the assumptions, in practice we may compare the patterns of estimated singular values under different lag values h_0 in iTOPUP and iTIPUP to evaluate the benefit of taking a larger h_0 . See also the simulation study.

We describe below the convergence rate of iTOPUP in terms of d_k , r_k and T under Assumption 3(i) when the dimensions of the core factor r_1, \dots, r_K are allowed to diverge.

COROLLARY 3.3. *Suppose Assumptions 1 and 3(i) hold. Let $h_0 \leq T/4$, $d_{-k}^* = \sum_{j \neq k} d_j r_j$ and $r = \prod_{k=1}^K r_k$. Suppose that for a sufficiently large C_0 not depending on $\{\sigma, d_k, r_k, k \leq K\}$,*

$$(3.25) \quad T \geq C_0 \max_{1 \leq k \leq K} \left(d^{2\delta_1 - \delta_0} r_k r_{-k}^2 + d^{2\delta_1} r_k^2 r_{-k} / d_k \right).$$

Then, after $J = O(\log d)$ iterations, we have the following upper bounds for iTOPUP,

$$(3.26) \quad \max_{1 \leq k \leq K} \|\widehat{P}_k^{(J)} - P_k\|_S = O_{\mathbb{P}}(1) \max_{1 \leq k \leq K} \left(\frac{d_k^{1/2} r^{3/2} (1 + r_k^{1/2} / d^{(1-\delta_0)/2})}{T^{1/2} d^{1/2 + \delta_0/2 - \delta_1} r_k} \right).$$

Moreover, (3.26) holds after at most $J = O(\log r)$ iterations, if any one of the following three conditions holds in addition to (3.25): (i) d_k ($k = 1, \dots, K$) are of the same order, (ii) λ_k ($k = 1, \dots, K$) are of the same order, (iii) $(\lambda_k)^{-2} \sqrt{d_k}$ ($k = 1, \dots, K$) are of the same order.

Note that the second part of Corollary 3.3 says that when the condition is right, iTOPUP algorithm only needs a small number of iterations to converge, as $O(\log r)$ is typically very small. The noise level σ does not appear directly in the rate since it is incorporated in the signal to noise ratio in the tensor form in Assumption 3. In Corollary 3.3, we show that as long as the sample size T satisfies (3.25), the iTOPUP achieves consistent estimation under regularity conditions. To digest (3.25), consider that the growth rate of r_k is much slower than d_k and the factors are strong with $\delta_0 = \delta_1 = 0$. Then (3.25) becomes $T \geq C_0 \max_k (r_k r_{-k}^2)$.

The advantage of using index δ_0, δ_1 is to link the convergence rates of the estimated factor loading space explicitly to the strength of factors. It is clear that the stronger the factors are, the faster the convergence rate is. Moreover, the stronger the factors are, the smaller the sample size is required.

When the ranks r_k ($k = 1, \dots, K$) also diverge and there is no severe signal cancellation in iTIPUP, we have the following convergence rate for iTIPUP under Assumption 3(ii).

COROLLARY 3.4. *Suppose Assumptions 1 and 3(ii) hold. Let $h_0 \leq T/4$ and $d_{-k}^* = \sum_{j \neq k} d_j r_j$. Suppose that for a sufficiently large C_0 not depending on $\{\sigma, d_k, r_k, k \leq K\}$,*

$$(3.27) \quad T \geq C_0 \max_{1 \leq k \leq K} \left(\frac{(d_k r_k + d^{\delta_0} r_k^2) r_{-k}^{2\delta_2} r^{2\delta_2}}{d^{1+3\delta_0-4\delta_1} \min_{1 \leq k \leq K} r_k^{2\delta_2}} + \frac{d_{-k}^* r_k r_{-k}^{2\delta_2}}{d^{1+\delta_0-2\delta_1}} \left(1 + \frac{r}{d^{1-\delta_0}}\right) \right).$$

Then, after at most $J = O(\log d)$ iterations, the iTIPUP estimator satisfies

$$(3.28) \quad \max_{1 \leq k \leq K} \|\hat{P}_k^{(J)} - P_k\|_S = O_{\mathbb{P}}(1) \max_{1 \leq k \leq K} \left(\frac{d_k^{1/2} r_k^{1/2} r_{-k}^{\delta_2} (1 + r^{1/2}/d^{(1-\delta_0)/2})}{T^{1/2} d^{1/2+\delta_0/2-\delta_1}} \right).$$

Moreover, (3.28) holds after at most $J = O(\log r)$ iterations, if any one of the following three conditions holds in addition to condition (3.27), (i) d_k ($k = 1, \dots, K$) are of the same order, (ii) λ_k^* ($k = 1, \dots, K$) are of the same order, (iii) $(\lambda_k^*)^{-2} \sqrt{d_k}$ ($k = 1, \dots, K$) are of the same order.

When the average auto-correlation is of unit order and the signal cancellation for TIPUP has no impact on the order of the signal ($\delta_0 = \delta_1$ and $\delta_2 = 0$ respectively), Corollary 3.4 requires the sampling rate $T \gtrsim h_0 + (d_k r_k + d^{\delta_0} r_k^2 + d_{-k}^* r_k (1 + r/d^{1-\delta_0}))/d^{1-\delta_0}$ and provides the convergence rate $(r_k d_k)^{1/2} (1 + r/d^{1-\delta_0})^{1/2} / (T d^{1-\delta_0})^{1/2}$. For examples, $T \geq 4h_0 + C_1$ gives the rate $(r_k d_k)^{1/2} / (T d^{1-\delta_0})^{1/2}$ when $\delta_0 \leq (K-1)/(2K)$ and $r_k^2 \lesssim d_k = d^{1/K} \forall k$, and the sample size requirement can be written as $T \gtrsim h_0 + d^{\delta_0} r_k^2 / d^{1-\delta_0}$ when $r_k^2 \asymp r^{2/K} \lesssim d_k \asymp d^{1/K} \forall k$ regardless of $\delta_0 \in [0, 1]$. Thus, the side condition involving $R^{*(\text{add})}$ in the second part of (3.17) is absorbed into the other components of (3.17).

Corollary 3.3 and Corollary 3.4 offer comparison of the iTOPUP and iTIPUP when the ranks diverge from two perspectives: sample size requirements and convergence rates. The lower bounds on T in (3.25) in Corollary 3.3 and (3.27) in Corollary 3.4 provide the sample complexity of the iTOPUP and iTIPUP respectively. In the case that the growth rate of r_k is much slower than d_k and the factors are strong with $\delta_0 = \delta_1 = 0$, the required sample size of the iTIPUP reduces to $T \geq 4h_0 + C_0 \max_{j,k} (r_k r_{-k}^{2\delta_2} r_{-j}^{2\delta_2} / d_{-k} + r_k r_{-k}^{2\delta_2} r_j / d_{-j})$, where $r_{-k} = r/r_k$ and $d_{-k} = d/d_k$. By comparing with the comment after Corollary 3.3, where the sample size requirement for the iTOPUP is $T \geq C_0 \max_k (r_k r_{-k}^2)$ when $\delta_0 = \delta_1 = 0$, it can be seen that the sample complexity for the iTIPUP is smaller, if δ_2 is a small constant. From the perspective of convergence rate, let us compare (3.26) in Corollary 3.3 and (3.28) in Corollary 3.4. When ranks diverge, iTIPUP is slower than iTOPUP if $\delta_2 > 3/2$ and faster if $\delta_2 \leq 1$, no matter how strong the factor is or what values δ_0, δ_1 take. As expected, the convergence rate is slower in the presence of weak factors. See the simulation for more empirical evidence.

Similar to Corollaries 3.3 and 3.4, we have the following rate for TIPUP-iTOPUP.

COROLLARY 3.5. *Suppose Assumptions 1 and 3 hold. Let $h_0 \leq T/4$ and $d_{-k}^* = \sum_{j \neq k} d_j r_j$. Suppose that for a sufficiently large C_0 not depending on $\{\sigma, d_k, r_k, k \leq K\}$,*

$$(3.29) \quad T \geq C_0 \max_{1 \leq k \leq K} \left(d^{2\delta_1 - \delta_0} r_k \left(\frac{r_{-k}^{2\delta_2}}{d_{-k}} + \frac{r_{-k}^3}{d_{-k}} \right) + \frac{d^{2\delta_1} r_k^2}{d_k} \left(\frac{r_{-k}^{2\delta_2}}{d_{-k}} + \frac{r_{-k}^3}{d_{-k}^2} \right) + \frac{d_{-k}^* \sqrt{r r_k}}{d^{1-\delta_1}} \right).$$

Then, after at most $J = O(\log d)$ iterations, the TIPUP-iTOPUP estimator satisfies

$$\max_{1 \leq k \leq K} \|\hat{P}_k^{(J)} - P_k\|_S = O_{\mathbb{P}}(1) \max_{1 \leq k \leq K} \left(\frac{d_k^{1/2} r^{3/2} (1 + r_k^{1/2}/d^{(1-\delta_0)/2})}{T^{1/2} d^{1/2+\delta_0/2-\delta_1} r_k} \right).$$

Moreover, the above error bound holds after at most $J = O(\log r)$ iterations, if any one of the following three conditions holds in addition to condition (3.29), (i) d_k ($k = 1, \dots, K$) are of the same order, (ii) λ_k ($k = 1, \dots, K$) are of the same order, (iii) $(\lambda_k)^{-2} \sqrt{d_k}$ ($k = 1, \dots, K$) are of the same order.

Compared with Corollary 3.3, Corollary 3.5 provides the same error bound for smaller T (possibly with bounded $T \gtrsim h_0$) when $r_{-k}^{2\delta_2} \lesssim r_{-k} d_{-k}$. The side condition involving $R^{(\text{add})}$ in the second part of (3.20), corresponding to the last component of (3.29) involving d_{-k}^* , is absorbed into the other components of (3.20) when $r_k^{1/2} \leq d^{\delta_1 - \delta_0} (r_{-k}^{2\delta_2 - 1} + r_{-k}^2) \forall k \leq K$.

3.5. Comparisons.

3.5.1. *Comparison between the non-iterative procedures and iterative procedures.* Theorems 3.1 and 3.2 show that the convergence rates of the non-iterative estimators TOPUP and TIPUP can be improved by their iterative counterparts. Particularly, when the dimensions r_k for the factor process are fixed and the respective signal strength conditions are fulfilled, the proposed iTOPUP and iTIPUP just need one-iteration to achieve the much sharper ideal rate $R^{(\text{ideal})}$ in (3.8) and $R^{*(\text{ideal})}$ in (3.15), compared with the rate (3.7) of TOPUP and (3.14) of TIPUP derived in Chen, Yang and Zhang (2019), respectively. The improvement is achieved through replacing the much larger d_{-k} by r_{-k} , via orthogonal projection. When the factors are strong with $\delta_0 = \delta_1 = 0$ and the factor dimensions are fixed, the non-iterative TOPUP-based estimators of Lam, Yao and Bathia (2011) for the vector factor model, Wang, Liu and Chen (2019) for the matrix factor and Chen, Yang and Zhang (2019) for tensor factor models all have the same $O_{\mathbb{P}}(T^{-1/2})$ convergence rate for estimating the loading space. In comparison, the convergence rate $O_{\mathbb{P}}(T^{-1/2} d_{-k}^{-1/2})$ of both iterative estimators, iTOPUP and iTIPUP (when there is no severe signal cancellation, with bounded δ_2), is much sharper. Intuitively, when the signal is strong, the orthogonal projection operation helps to consolidate signals while potentially averaging out the noises, when the projection reduces the dimension of the mode- k unfolded matrix from $d_k \times d_{-k}$ for the tensor \mathcal{X}_t to $d_k \times r_{-k}$ for the projected tensor \mathcal{Z}_t , resulting in the improvement by a factor of $d_{-k}^{-1/2}$ in the convergence rate.

When r_k are allowed to diverge, the iTOTUP and iTIPUP algorithms converge after at most $O(\log(d))$ iterations to achieve the ideal rate according to Theorems 3.1 and 3.2. The number of iterations needed can be as few as $O(\log(r))$ when the condition is right.

3.5.2. *Comparison between iTIPUP and iTOPUP.* The inner product operation in (2.7) for TIPUP-related procedures enjoys significant amount of noise cancellation comparing to the outer product operation in (2.4) for TOPUP-related procedures. Compared with iTOPUP, the benefit of noise cancellation of the iTIPUP procedure is still visible through the reduction of r_{-k} in (3.8) to $\sqrt{r_{-k}}$ in (3.15) in the ideal rates. However, this post-iteration benefit is much less pronounced compared with the reduction of d_{-k} in (3.7) for TOPUP to $\sqrt{d_{-k}}$ in (3.14) for TIPUP in the non-iterative rates. Meanwhile, the potential for signal cancellation in the TIPUP related schemes persists as λ_k^* and λ_k are unchanged between the initial and ideal rates. We note that the signal strength can be viewed as λ_k and λ_k^* in Theorems 3.1 and 3.2 respectively for TOPUP/iTOPUP and TIPUP/iTIPUP, and that severe signal cancellation can be expressed as $\lambda_k^* \ll \lambda_k$. When r_{-k} are allowed to diverge to infinity, the impact of signal cancellation is expressed in terms of δ_2 in Assumption 3: The iTOPUP has a faster rate than the iTIPUP when $\delta_2 > 3/2$ and slower rate when $\delta_2 \leq 1$, in view of Corollary 3.3 and 3.4. In Corollaries 3.1 and 3.2, iTOPUP and iTIPUP have the same convergence rate because Corollary 3.2 assumes that signal cancellation does not change convergence rate.

Our results seem to suggest that the mixed TIPUP-iTOPUP procedure would strike a good balance between the benefit of noise cancellation (e.g. smaller T for consistency) and the potential danger of signal cancellation (e.g. $\lambda_k^* \ll \lambda_k$) for the following four reasons: (1) The benefit of noise cancellation is much larger in the initialization, in term of d_{-k} , in view of the rates $R_k^{(0)}$ in (3.7) and $R^{*(0)}$ in (3.14). (2) The first part of condition (3.20) for TIPUP-iTOPUP is weaker than the first part of condition (3.17) for TIPUP-iTIPUP. (3) The signal strength λ_k of the stronger TOPUP form is retained in the rate $R^{(\text{ideal})}$ after iTOPUP iteration. (4) As we will prove in Section 3.6, the sample size requirement for the TIPUP initialization is optimal in the sense that it matches a computational lower bound under suitable conditions. Our simulation results support this recommendation, especially for relatively small r_{-k} . Of course if the sample size qualitatively justifies the condition $C_1^{(\text{TOPUP})} R^{(0)} \leq (1 - \rho)/4$ in (3.10) and/or if a possible signal cancellation is a significant concern, the TOPUP initiation should be used.

3.5.3. Comparison with HOOI. The signal to noise ratio (SNR) condition, or equivalently the sample size requirement, is mainly used to ensure that the initial estimator has sufficiently small estimation error. Thus, the performance of iterative procedures is measured by both the SNR requirement and the error rate achieved. Consider fixed h_0 in the fixed rank case with $K = 3$ and $d_{\max} \asymp d^{1/K}$. In the fixed signal model where $\mathcal{M}_t = \mathcal{M}$ is fixed and deterministic in (1.1), applying HOOI to the average of \mathcal{X}_t would require SNR $\lambda(T^{1/2}/\sigma) \geq C_0 d^{1/4}$ to achieve the loss of the order $(\sigma/T^{1/2})d_k^{1/2}/\lambda$ according to Zhang and Xia (2018), where $\sigma/T^{1/2}$ is viewed as the noise level for HOOI as it is the standard deviation of each element of the average tensor. In terms of the auto-crossproducts, taking the average over \mathcal{X}_t roughly amounts to taking the average of all $T(T-1)/2$ lagged products between \mathcal{X}_{t-h} and \mathcal{X}_t , $1 \leq t-h < t \leq T$. However, in the tensor factor model (1.1) where the signal part is random and serial correlated, the average is taken only over $T-h$ lagged products for each h . Thus, while the rate of the average of the signal-by-noise crossproducts in the factor model is heuristically expected to match that of HOOI at noise level $\sigma/T^{1/2}$, the rate of the average of the noise-by-noise crossproducts in the factor model is expected to only match that of HOOI with noise level $\sigma/T^{1/4}$. In Corollary 3.2, the contribution of the noise-by-noise crossproducts dominates the initial estimation error as the SNR requirement $\lambda(T^{1/4}/\sigma) \geq C_0 d^{1/4}$ in (3.23) matches that of HOOI with noise level $\sigma/T^{1/4}$; at the same time the contribution of the signal-by-noise crossproducts dominates the estimation error after iteration as the rate $(\sigma/T^{1/2})d_k^{1/2}/\lambda$ in (3.24) matches that of HOOI with noise level $\sigma/T^{1/2}$. Thus, if there is no severe signal cancellation, the signal to noise ratio requirement and convergence rate for iTIPUP and TIIP-iTOPUP in the factor model are both comparable with those of HOOI in the simpler fixed signal setting, but the rate match is achieved in very different and subtle ways. We prove that this insight is intrinsic as the rates in (3.23) and (3.24) are both optimal according to the computational and statistical lower bounds in the following subsection.

3.6. Computational and statistical lower bounds. In this subsection, we focus on the typical factor model setting that the condition numbers of $A_k^\top A_k$ are bounded and ranks r_k are fixed. We shall prove that under the computational hardness assumption, the signal to noise ratio condition (3.23) imposed on iTIPUP (also TIPUP-iTOPUP) in Corollary 3.2 is unavoidable for computationally feasible estimators to be consistent. To be specific, we show that, if the signal to noise ratio condition is violated, then any computationally efficient and consistent estimator of the loading spaces leads to a computationally efficient and statistically consistent test for the Hypergraphic Planted Clique Detection problem in a regime where it is

believed to be computationally intractable. In addition, we establish a statistical lower bound on the minimax risk of the estimators.

Hypergraphic Planted Clique. An m -hypergraph $G = (V(G), E(G))$ is a natural extension of regular graph, where $V(G) = [N]$ and each hyper-edge is represented by an unordered group of m different vertices $i_j \in V(G)$ ($j = 1, \dots, m$), denoted as $e = (i_1, \dots, i_m) \in E(G)$. Given a m -hypergraph its adjacency tensor $\mathcal{A} \in \{0, 1\}^{N \times N \times \dots \times N}$ is defined as

$$\mathcal{A}_{i_1, \dots, i_m} = \begin{cases} 1, & \text{if } e = (i_1, \dots, i_m) \in E(G); \\ 0, & \text{otherwise.} \end{cases}$$

We denote by $\mathcal{G}_m(N, 1/2)$ the Erdős–Rényi m -hypergraph on N vertices where each hyper-edge e is drawn independently with probability $1/2$, by $\mathcal{C} = \mathcal{C}(N, \kappa)$ a random clique of size κ where the κ members are uniformly sampled from $[N]$ and $E(\mathcal{C})$ is composed of all $e = (i_1, \dots, i_m)$ with $i_j \in \mathcal{C}$, and by $\mathcal{G}_m(N, 1/2, \kappa)$ the random graph generated by first sampling independently $\mathcal{G}_m(N, 1/2)$ and $\mathcal{C} = \mathcal{C}(N, \kappa)$ and then adding all the edges in $E(\mathcal{C})$ to the set of edges in $\mathcal{G}_m(N, 1/2)$. The Hypergraphic Planted Clique (HPC) detection problem of parameter (N, κ, m) refers to testing the following hypotheses:

$$(3.30) \quad H_0^G : \mathcal{A} \sim \mathcal{G}_m(N, 1/2) \quad \text{v.s.} \quad H_1^G : \mathcal{A} \sim \mathcal{G}_m(N, 1/2, \kappa).$$

If $m = 2$, the above HPC detection becomes the traditional planted clique (PC) detection problem. When $\kappa \geq c\sqrt{N}$, many computationally efficient algorithms have been developed for PC detection; see, [Alon, Krivelevich and Sudakov \(1998\)](#), [Feige and Krauthgamer \(2000\)](#), [Feige and Ron \(2010\)](#), [Ames and Vavasis \(2011\)](#), [Dekel, Gurel-Gurevich and Peres \(2014\)](#), [Deshpande and Montanari \(2015\)](#), [Feldman et al. \(2017\)](#), among others. However, it has been widely conjectured that when $\kappa = o(\sqrt{N})$, the PC detection problem cannot be solved in randomized polynomial time, which is referred to as the hardness conjecture. Computational lower bounds in several statistical problems have been established by assuming the hardness conjecture of PC detection, including sparse PCA ([Berthet and Rigollet, 2013a,b](#), [Wang, Berthet and Samworth, 2016](#)), sparse CCA ([Gao, Ma and Zhou, 2017](#)), submatrix detection ([Ma and Wu, 2015](#), [Cai, Liang and Rakhlin, 2017](#)), community detection ([Hajek, Wu and Xu, 2015](#)), etc.

Recently, motivated by tensor data analysis, hardness conjecture for HPC detection problem has been proposed; see, for example, [Zhang and Xia \(2018\)](#), [Brennan and Bresler \(2020\)](#), [Luo and Zhang \(2020a,b\)](#), [Pananjady and Samworth \(2020\)](#). Similar to the PC detection, they hypothesized that when $\kappa = O(N^{1/2-\delta})$ with $\delta > 0$, the HPC detection problem (3.30) cannot be solved by any randomized polynomial-time algorithm. Formally, the conjectured hardness of the HPC detection problem can be stated as follows.

HYPOTHESIS I (HPC detection). Consider the HPC detection problem (3.30) and suppose $m \geq 2$ is a fixed integer. If

$$(3.31) \quad \limsup_{N \rightarrow \infty} \frac{\log \kappa}{\log N} \leq \frac{1}{2} - \delta, \quad \text{for any } \delta > 0,$$

for any sequence of polynomial-time tests $\{\psi\}_N : \mathcal{A} \rightarrow \{0, 1\}$,

$$\limsup_{N \rightarrow \infty} (\mathbb{P}_{H_0^G}(\psi(\mathcal{A}) = 1) + \mathbb{P}_{H_1^G}(\psi(\mathcal{A}) = 0)) > 1/2.$$

Evidence supporting this hypothesis has been provided in [Zhang and Xia \(2018\)](#), [Luo and Zhang \(2020a\)](#). This version of the hypothesis is similar to the one in [Berthet and Rigollet \(2013a\)](#), [Ma and Wu \(2015\)](#), [Gao, Ma and Zhou \(2017\)](#) for the PC detection problem.

For simplicity, we especially consider the one-factor model (1.2) with \mathcal{F}_t being a mean 0 univariate series,

$$(3.32) \quad \mathcal{X}_t = \lambda f_t \times_1 a_1 \times_2 \dots \times_K a_K + \mathcal{E}_t,$$

where $a_k \in \mathbb{R}^{d_k}$, $\|a_k\|_2 = 1$ for $1 \leq k \leq K$, and $\mathbb{E}f_t^2 = 1$. The probability space we consider in this subsection is

$$(3.33) \quad \mathcal{P}(T, d_1, \dots, d_K, \lambda) = \left\{ \mathcal{X}_1, \dots, \mathcal{X}_T : \text{each } \mathcal{X}_t \text{ has form (3.32) with } f_t \sim N(0, 1), \right. \\ \left. \frac{1}{T-1} \sum_{t=2}^T \mathbb{E}f_t f_{t-1} = c_0 > 0, \text{ and } \{f_t\}_{t=1}^T \text{ independent of } \{\mathcal{E}_t\}_{t=1}^T, \right. \\ \left. \mathcal{E}_{t, j_1, \dots, j_K} \stackrel{i.i.d.}{\sim} N(0, 1) \text{ for all } 1 \leq t \leq T, 1 \leq j_k \leq d_k, 1 \leq k \leq K \right\}.$$

The computational lower bound over $\mathcal{P}(T, d_1, \dots, d_K, \lambda)$ is then presented as below.

THEOREM 3.4. *Suppose that Hypothesis 1 holds for some $0 < \delta < 1/2$ and $d^{1/K} \asymp d_k \geq T$ for all $1 \leq k \leq K$. If, for some $\vartheta > 0$,*

$$(3.34) \quad \liminf_{T \rightarrow \infty} \frac{\sigma^2 d^{1/2-\vartheta}}{T^{1/2} \lambda^2} > 0,$$

then for any randomized polynomial-time estimators $\hat{a}_k = \hat{a}_k(\mathcal{X}_1, \dots, \mathcal{X}_T)$, $1 \leq k \leq K$,

$$(3.35) \quad \liminf_{T \rightarrow \infty} \sup_{\mathcal{X}_1, \dots, \mathcal{X}_T \in \mathcal{P}(T, d_1, \dots, d_K, \lambda)} \mathbb{P} \left(\min_{1 \leq k \leq K} \|\hat{P}_k - P_k\|_S^2 > \frac{1}{3} \right) > \frac{1}{4},$$

where $\hat{P}_k = \hat{a}_k \hat{a}_k^\top$ and $P_k = a_k a_k^\top$.

Comparing (3.34) with (3.23), we see that the signal to noise ratio condition (3.23) cannot be improved upon by a factor of d^ϑ with polynomial time complexity for any $\vartheta > 0$.

REMARK 3.9. Theorem 3.4 illustrates the computational hardness for factor loading spaces estimation under the typical factor model setting that the condition numbers of $A_k^\top A_k$ are bounded and ranks r_k are fixed, and suggests the use of TIPUP initialization with proper fixed h_0 as it attains the computational lower bound under the typical factor model setting.

Next, we establish the statistical lower bound for the tensor factor model problem. Again, we consider the probability space (3.33).

THEOREM 3.5. *Suppose $\lambda > 0$ and $d_k \rightarrow \infty$ as $T \rightarrow \infty$ for all $1 \leq k \leq K$. Then there exists a universal constant $c > 0$ such that for T sufficiently large,*

$$(3.36) \quad \inf_{\hat{a}_k} \sup_{\mathcal{X}_1, \dots, \mathcal{X}_T \in \mathcal{P}(T, d_1, \dots, d_K, \lambda)} \mathbb{E} \|\hat{P}_k - P_k\|_S \geq c \min \left(1, (\sigma^2 + \sigma \lambda) \sqrt{d_k} / (\lambda^2 \sqrt{T}) \right)$$

for all $1 \leq k \leq K$, where $\hat{P}_k = \hat{a}_k \hat{a}_k^\top$ and $P_k = a_k a_k^\top$.

REMARK 3.10. Theorem 3.5 provides statistical lower bound for high dimensional tensor factor models. It matches the upper bounds in Corollary 3.1 and 3.2, showing that the rates obtained by our proposed iterative procedures are minimax-optimal.

4. Summary. In this paper we propose new estimation procedures for tensor factor model via iterative projection, and focus on two procedures: iTOPUP and iTIPUP. Theoretical analysis shows the asymptotic properties of the estimators. Simulation study presented in the supplementary material illustrates the finite sample properties of the estimators. While theoretical results are obtained under very general conditions, concrete specific cases are considered. In particular, under the typical factor model setting where the condition numbers of $A_k^\top A_k$ are bounded and the ranks r_k are fixed, the proposed iterative procedures, iTOPUP method and iTIPUP method (with no severe signal cancellation) lead to a convergence rate $O_{\mathbb{P}}((Td_{-k})^{-1/2})$ under strong factors settings due to information pooling of the orthogonal projection of the other d_{-k} dimensions. This rate is much sharper than the existing rate $O_{\mathbb{P}}(T^{-1/2})$ in the recent literature for non-iterative estimators for vector, matrix and tensor factor models. It implies that the accuracy can be improved by increasing the dimensions, and consistent estimation of the loading spaces can be achieved even with a fixed finite sample size T . This is in sharp contrast to the folklore based on the existing literature that only the sample size T helps the estimation of the loading matrices in factor models. The proposed iterative estimation methods not only preserve the tensor structure, but also result in sharper convergence rate in the estimation of factor loading space.

The iterative procedure requires two operators, one for initialization and one for iteration. Under certain conditions of the signal to noise ratio (or the sample size requirement), we only need the initial estimator to have sufficiently small estimation errors but not the consistency of the initial estimator. Often, one iteration is sufficient. In more complicated general cases, at most $O(\log(d))$ iterations are needed to achieve the ideal rate of convergence. Based on the theoretical results and empirical evidence, we suggest to use iTOPUP for iteration when the ranks r_k are small. In terms of initiation, the computational lower bound shows that the signal to noise ratio condition derived from TIPUP initialization is unavoidable for any computationally feasible estimation procedure to achieve consistency, while that from TOPUP initialization is not optimal. Based on this result, we suggest the use of TIPUP initialization. Of course, this should be done with precaution against potential signal cancellation, for example by using a slightly large h_0 as our empirical results show. By examination of the patterns of estimated singular values under different lag values h_0 , using iTOPUP and iTIPUP, it is possible to detect signal cancellation, which has significant impact on iTIPUP estimators.

The proposed iterative procedure is similar to HOOI algorithms in spirit, but the detailed operations and the theoretical challenges are significantly different.

Acknowledgements. We would like to thank the Editor, the Associate Editor and the anonymous referees for their detailed reviews, which helped to improve the paper substantially.

Yuefeng Han’s research is supported in part by National Science Foundation grant IIS-1741390. Rong Chen’s research is supported in part by National Science Foundation grants DMS-1737857, IIS-1741390, CCF-1934924 and DMS-2027855. Dan Yang’s research is supported in part by NSF grant IIS-1741390, Hong Kong grant GRF 17301620 and Hong Kong grant CRF C7162-20GF. Cun-Hui Zhang’s research is supported in part by NSF grants DMS-1721495, IIS-1741390 and CCF-1934924.

REFERENCES

AHN, S. C. and HORENSTEIN, A. R. (2013). Eigenvalue ratio test for the number of factors. *Econometrica* **81** 1203–1227.
 ALON, N., KRIVELEVICH, M. and SUDAKOV, B. (1998). Finding a large hidden clique in a random graph. *Random Structures & Algorithms* **13** 457–466.

- ALTER, O. and GOLUB, G. H. (2005). Reconstructing the pathways of a cellular system from genome-scale signals by using matrix and tensor computations. *Proceedings of the National Academy of Sciences* **102** 17559–17564.
- AMES, B. P. and VAVASIS, S. A. (2011). Nuclear norm minimization for the planted clique and biclique problems. *Mathematical Programming* **129** 69–89.
- ANANDKUMAR, A., GE, R., HSU, D. and KAKADE, S. M. (2014). A tensor approach to learning mixed membership community models. *The Journal of Machine Learning Research* **15** 2239–2312.
- BAI, J. (2003). Inferential theory for factor models of large dimensions. *Econometrica* **71** 135–171.
- BAI, J. and NG, S. (2002). Determining the number of factors in approximate factor models. *Econometrica* **70** 191–221.
- BAI, J. and NG, S. (2007). Determining the number of primitive shocks in factor models. *Journal of Business & Economic Statistics* **25** 52–60.
- BERTHET, Q. and RIGOLLET, P. (2013a). Complexity theoretic lower bounds for sparse principal component detection. In *Conference on Learning Theory* 1046–1066. PMLR.
- BERTHET, Q. and RIGOLLET, P. (2013b). Optimal detection of sparse principal components in high dimension. *Annals of Statistics* **41** 1780–1815.
- BI, X., QU, A. and SHEN, X. (2018). Multilayer tensor factorization with applications to recommender systems. *Annals of Statistics* **46** 3308–3333.
- BIRNBAUM, A., JOHNSTONE, I. M., NADLER, B. and PAUL, D. (2013). Minimax bounds for sparse PCA with noisy high-dimensional data. *Annals of Statistics* **41** 1055.
- BRENNAN, M. and BRESLER, G. (2020). Reducibility and statistical-computational gaps from secret leakage. In *Conference on Learning Theory* 648–847. PMLR.
- CAI, T. T., LIANG, T. and RAKHLIN, A. (2017). Computational and statistical boundaries for submatrix localization in a large noisy matrix. *Annals of Statistics* **45** 1403–1430.
- CHAMBERLAIN, G. and ROTHSCILD, M. (1983). Arbitrage, Factor Structure, and Mean-Variance Analysis on Large Asset Markets. *Econometrica* **51** 1281–1304.
- CHEN, E. Y. and CHEN, R. (2019). Modeling Dynamic Transport Network with Matrix Factor Models: with an Application to International Trade Flow. *arXiv preprint arXiv:1901.00769*.
- CHEN, E. Y., FAN, J. and LI, E. (2020). Statistical Inference for High-Dimensional Matrix-Variate Factor Model. *arXiv preprint arXiv:2001.01890*.
- CHEN, E. Y., TSAY, R. S. and CHEN, R. (2019). Constrained Factor Models for High-Dimensional Matrix-Variate Time Series. *Journal of the American Statistical Association* 1–37.
- CHEN, R., YANG, D. and ZHANG, C.-H. (2019). Factor Models for High-Dimensional Tensor Time Series. *arXiv preprint arXiv:1905.07530*.
- CHEN, E. Y., XIA, D., CAI, C. and FAN, J. (2020). Semiparametric tensor factor analysis by iteratively projected svd. *arXiv preprint arXiv:2007.02404*.
- DE LATHAUWER, L., DE MOOR, B. and VANDEWALLE, J. (2000). On the best rank-1 and rank-(r_1, r_2, \dots, r_n) approximation of higher-order tensors. *SIAM Journal on Matrix Analysis and Applications* **21** 1324–1342.
- DEKEL, Y., GUREL-GUREVICH, O. and PERES, Y. (2014). Finding hidden cliques in linear time with high probability. *Combinatorics, Probability and Computing* **23** 29–49.
- DESHPANDE, Y. and MONTANARI, A. (2015). Finding hidden cliques of size $\sqrt{N/e}$ in nearly linear time. *Foundations of Computational Mathematics* **15** 1069–1128.
- DIACONIS, P. and FREEDMAN, D. (1980). Finite exchangeable sequences. *Annals of Probability* 745–764.
- FAN, J., LIAO, Y. and MINCHEVA, M. (2011). High dimensional covariance matrix estimation in approximate factor models. *Annals of Statistics* **39** 3320.
- FAN, J., LIAO, Y. and MINCHEVA, M. (2013). Large covariance estimation by thresholding principal orthogonal complements. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **75** 603–680.
- FAN, J., LIAO, Y. and WANG, W. (2016). Projected principal component analysis in factor models. *Annals of Statistics* **44** 219–254.
- FAN, J., LIU, H. and WANG, W. (2018). Large covariance estimation through elliptical factor models. *Annals of Statistics* **46** 1383.
- FEIGE, U. and KRAUTHGAMER, R. (2000). Finding and certifying a large hidden clique in a semirandom graph. *Random Structures & Algorithms* **16** 195–208.
- FEIGE, U. and RON, D. (2010). Finding hidden cliques in linear time. In *Discrete Mathematics and Theoretical Computer Science* 189–204. Discrete Mathematics and Theoretical Computer Science.
- FELDMAN, V., GRIGORESCU, E., REYZIN, L., VEMPALA, S. S. and XIAO, Y. (2017). Statistical algorithms and a lower bound for detecting planted cliques. *Journal of the ACM (JACM)* **64** 1–37.
- FORNI, M., HALLIN, M., LIPPI, M. and REICHLIN, L. (2000). The Generalized Dynamic-Factor Model: Identification and Estimation. *The Review of Economics and Statistics* **82** 540–554.

- FOSTER, G. (1996). Time series analysis by projection. II. Tensor methods for time series analysis. *The Astrophysical Journal* **111** 555.
- GAO, C., MA, Z. and ZHOU, H. H. (2017). Sparse CCA: Adaptive estimation and computational barriers. *Annals of Statistics* **45** 2074–2101.
- HAJEK, B., WU, Y. and XU, J. (2015). Computational lower bounds for community detection on random graphs. In *Conference on Learning Theory* 899–928. PMLR.
- HALLIN, M. and LIŠKA, R. (2007). Determining the number of factors in the general dynamic factor model. *Journal of the American Statistical Association* **102** 603–617.
- LAM, C., YAO, Q. and BATHIA, N. (2011). Estimation of latent factors for high-dimensional time series. *Biometrika* **98** 901–918.
- LAM, C. and YAO, Q. (2012). Factor modeling for high-dimensional time series: inference for the number of factors. *Annals of Statistics* **40** 694–726.
- LIU, J., MUSIALSKI, P., WONKA, P. and YE, J. (2012). Tensor completion for estimating missing values in visual data. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **35** 208–220.
- LIU, Y., SHANG, F., FAN, W., CHENG, J. and CHENG, H. (2014). Generalized higher-order orthogonal iteration for tensor decomposition and completion. In *Advances in Neural Information Processing Systems* 1763–1771.
- LUO, Y. and ZHANG, A. R. (2020a). Tensor clustering with planted structures: Statistical optimality and computational limits. *arXiv preprint arXiv:2005.10743*.
- LUO, Y. and ZHANG, A. R. (2020b). Open problem: Average-case hardness of hypergraphic planted clique detection. In *Conference on Learning Theory* 3852–3856. PMLR.
- MA, Z. and WU, Y. (2015). Computational barriers in minimax submatrix detection. *Annals of Statistics* **43** 1089–1116.
- OMBERG, L., GOLUB, G. H. and ALTER, O. (2007). A tensor higher-order singular value decomposition for integrative analysis of DNA microarray data from different studies. *Proceedings of the National Academy of Sciences* **104** 18371–18376.
- PAN, J. and YAO, Q. (2008). Modelling multiple time series via common factors. *Biometrika* **95** 365–379.
- PANANJADY, A. and SAMWORTH, R. J. (2020). Isotonic regression with unknown permutations: Statistics, computation, and adaptation. *arXiv preprint arXiv:2009.02609*.
- PENA, D. and BOX, G. E. (1987). Identifying a simplifying structure in time series. *Journal of the American Statistical Association* **82** 836–843.
- ROGERS, M., LI, L. and RUSSELL, S. J. (2013). Multilinear dynamical systems for tensor time series. *Advances in Neural Information Processing Systems* **26** 2634–2642.
- SHEEHAN, B. N. and SAAD, Y. (2007). Higher order orthogonal iteration of tensors (HOOI) and its relation to PCA and GLRAM. In *Proceedings of the 2007 SIAM International Conference on Data Mining* 355–365. SIAM.
- STOCK, J. H. and WATSON, M. W. (2002). Forecasting Using Principal Components from a Large Number of Predictors. *Journal of the American Statistical Association* **97** 1167–1179.
- SUN, W. W. and LI, L. (2017). STORE: sparse tensor response regression and neuroimaging analysis. *The Journal of Machine Learning Research* **18** 4908–4944.
- TUCKER, L. R. (1966). Some mathematical notes on three-mode factor analysis. *Psychometrika* **31** 279–311.
- WANG, T., BERTHET, Q. and SAMWORTH, R. J. (2016). Statistical and computational trade-offs in estimation of sparse principal components. *Annals of Statistics* **44** 1896–1930.
- WANG, D., LIU, X. and CHEN, R. (2019). Factor models for matrix-valued high-dimensional time series. *Journal of Econometrics* **208** 231–248.
- WANG, D., ZHENG, Y. and LI, G. (2021). High-Dimensional Low-Rank Tensor Autoregressive Time Series Modeling. *arXiv preprint arXiv:2101.04276*.
- WEDIN, P.-Å. (1972). Perturbation bounds in connection with singular value decomposition. *BIT Numerical Mathematics* **12** 99–111.
- ZHANG, A. and XIA, D. (2018). Tensor svd: Statistical and computational limits. *IEEE Transactions on Information Theory* **64** 7311–7338.
- ZHOU, H., LI, L. and ZHU, H. (2013). Tensor regression with applications in neuroimaging data analysis. *Journal of the American Statistical Association* **108** 540–552.

**SUPPLEMENTARY MATERIAL TO “TENSOR FACTOR MODEL
ESTIMATION BY ITERATIVE PROJECTION”**

BY YUEFENG HAN, RONG CHEN, DAN YANG, AND CUN-HUI ZHANG

Rutgers University and The University of Hong Kong

In this supplementary material, we shall provide simulation studies, the proofs of main results in the paper and some lemmas that are useful in proofs of the paper.

The readers are referred to Appendix A for simulation studies. The proofs of Theorems 3.2, 3.1, 3.4 and 3.5 are presented in Appendix B, C, D and E, respectively. Appendix F includes the proofs of corollaries. All technical lemmas are relegated to Appendix G.

APPENDIX A: SIMULATION STUDY

In this section, we compare the empirical performance of different procedures of estimating the loading matrices of a tensor factor model, under various simulation setups. Specifically, we consider the following procedures: the non-iterative and iterative methods, and the intermediate output from the iterative procedures when the number of iteration is 1 after initialization. If TIPUP is used as UINIT and UITER, the one step procedure will be denoted as 1TIPUP. Similarly for 1UP and 1TOPUP. We consider the following combinations of UINIT and UITER.

- UP based: (i) UP, (ii) 1UP and (iii) iUP
- TIPUP based: (iv) TIPUP, (v) 1TIPUP and (vi) iTIPUP
- TOPUP based: (vii) TOPUP, (viii) 1TOPUP and (ix) iTOPUP
- mixed iterative: (x) TIPUP-1TOPUP, and (xi) TIPUP-iTOPUP
- mixed iterative: (xii) TOPUP-1TIPUP, and (xiii) TOPUP-iTIPUP

Our empirical results show that (xii) and (xiii) are always inferior to (v) and (vi) respectively. Hence results used (xii) and (xiii) are not shown here.

We demonstrate the performance of all procedures under the setting of a matrix factor model,

$$(A.1) \quad X_t = A_1 F_t A_2^\top + E_t = \lambda U_1 F_t U_2^\top + E_t,$$

and an order 3 tensor factor model

$$(A.2) \quad \mathcal{X}_t = \lambda \mathcal{F}_t \times_1 U_1 \times_2 U_2 \times_3 U_3 + \mathcal{E}_t.$$

Under matrix factor model (A.1), we consider two experimental configurations:

- I. We set $r_1 = r_2 = 1$ for three purposes: to see the effect of sample size T and signal strength λ ; to check the effect of h_0 when there is no signal cancellation; and to verify the theoretical bounds on the sample size.
- II. We consider the rank setup $r_1 = 1$ and $r_2 = 2$ while fixing the signal strength λ . We vary the parameters so that signal cancellation may or may not happen, and under the latter case, the choice of h_0 plays an important role.

Under tensor factor model (A.2), we implement the following configuration:

- III. The ranks are given as $r_1 = r_2 = r_3 = 2$.

We repeat all the experiments 100 times.

Here, E_t in matrix factor model (A.1) (resp. \mathcal{E}_t in tensor factor model (A.2)) is white noise with no autocorrelation, $E_t \perp E_{t+h}$, $h > 0$ (resp. $\mathcal{E}_t \perp \mathcal{E}_{t+h}$, $h > 0$), and generated according

to $E_t = \Psi_1^{1/2} Z_t \Psi_2^{1/2}$ (resp. $\mathcal{E}_t = Z_t \times_1 \Psi_1^{1/2} \times_2 \Psi_2^{1/2} \times_3 \Psi_3^{1/2}$), where all of the elements in the $d_1 \times d_2$ matrix Z_t (resp. $d_1 \times d_2 \times d_3$ tensor Z_t) are iid $N(0, 1)$. Furthermore, Ψ_1, Ψ_2 (resp. Ψ_3) are the covariance matrices along each mode with the diagonal elements being 1 and all the off-diagonal elements being ψ_1, ψ_2 (resp. ψ_3). The elements of the loading matrices U_j of size $d_j \times r_j$, for $j = 1, 2, 3$, are first generated from iid $N(0, 1)$, and then orthonormalized through QR decomposition. We set different λ values for different signal-to-noise ratio.

As for the factor time series, under matrix factor model Configuration I, the univariate f_t follows AR(1) with AR coefficient ϕ ; under matrix factor model Configuration II, F_t is a 1×2 matrix and the two univariate time series f_{it} follow AR(1) $f_{it} = \phi_i f_{i(t-1)} + \epsilon_{it}$ independently with two AR coefficients ϕ_1 and ϕ_2 respectively; under tensor factor model Configuration III, \mathcal{F}_t is a $2 \times 2 \times 2$ tensor with eight independent univariate time series, where three follow AR(1) processes $f_{111t} = 0.7f_{111t-1} + e_{111t}$, $f_{211t} = 0.6f_{211t-1} + e_{211t}$, $f_{222t} = 0.8f_{222t-1} + e_{222t}$, one follows AR(2) process $f_{221t} = 0.5f_{221t-1} + 0.3f_{221t-2} + e_{221t}$, and four, $f_{112t}, f_{121t}, f_{122t}, f_{212t}$, are white noise. Here, all of the innovations follow iid $N(0, 1)$.

We fix dimensions $d_1 = d_2 (= d_3) = 16$ and use the estimation error for A_1 or U_1 as criterion: $\|\hat{P}_1 - P_1\|_S$. The λ in (A.1) and (A.2) is $\lambda = \prod_{k=1}^K \|A_k\|_S$ in Corollaries 3.1 and 3.2.

Configuration I satisfies Assumption 2 since the rank r_1 and r_2 are fixed and the factor process is stationary. We performed three experiments under Configuration I for three purposes.

Experiment 1 under Configuration I. We set the off-diagonal entries of the covariance matrices of the noise as $\psi_1 = \psi_2 = 0.2$ and the AR coefficient of the factor time series as $\phi = 0.8$, and vary the sample size $T = 256, 1024$ and signal strength $\lambda = 1, 2, 4$. Figure 1 shows the boxplot of the logarithm of the estimation errors for methods (i)-(ix). The performance of the mixed algorithms (x)-(xi) is not shown because they are identical to the corresponding methods (v) and (vi), under the rank one setting when the same initialization is used. We use $h_0 = 1$ and $\hat{r}_2 = 1$ in the process of the estimation. It can be seen easily from Figure 1 that UP, 1UP, and iUP are always the worst, showing the advantage of the methods that accommodate time series features and the disadvantage of neglecting temporal correlation. We will exclude UP, 1UP, and iUP from comparison for the rest of the simulation. When the sample size is small and the signal is weak ($T = 256$ and $\lambda = 1$), none of the methods work well, though procedures using TIPUP work sometimes. When the sample size is not too small or the signal strength is not too weak (shown in all panels except for the top left one), one-step methods (1TIPUP and 1TOPUP) are better than the noniterative methods (TIPUP and TOPUP), and iterative methods (iTIPUP and iTOPUP) are in turn better than the one-step methods. When the sample size and signal strength increase, all methods perform better, but meanwhile the advantage of iterative methods over one-step methods and the advantage of one-step methods over initialization methods become smaller. When the sample size is large and signal is strong, the one-step methods are similar to the iterative methods after convergence, corroborating Corollaries 3.1 and 3.2.

It is somewhat surprising to observe, from the top left panel in Figure 1 ($T = 256$ and $\lambda = 1$), that in the small sample size and low signal strength case, the median error of iTIPUP is larger than that of 1TIPUP, which in turn is larger than TIPUP, whereas the order is reversed under stronger signal to noise ratio or with larger sample size shown in the other panels. Furthermore, the top right panel in Figure 1 shows that, with weak signal to noise ratio, the TIPUP based methods perform better than the TOPUP based methods. This observation coincides with the results in Corollaries 3.1 and 3.2, which together state that iTOPUP requires larger signal-to-noise ratio for consistency than iTIPUP. Figure 2 produces some deeper insight, where the trajectories of the iterative methods (including initial estimations, estimations after one iteration, and the estimations after final convergence) of the 100

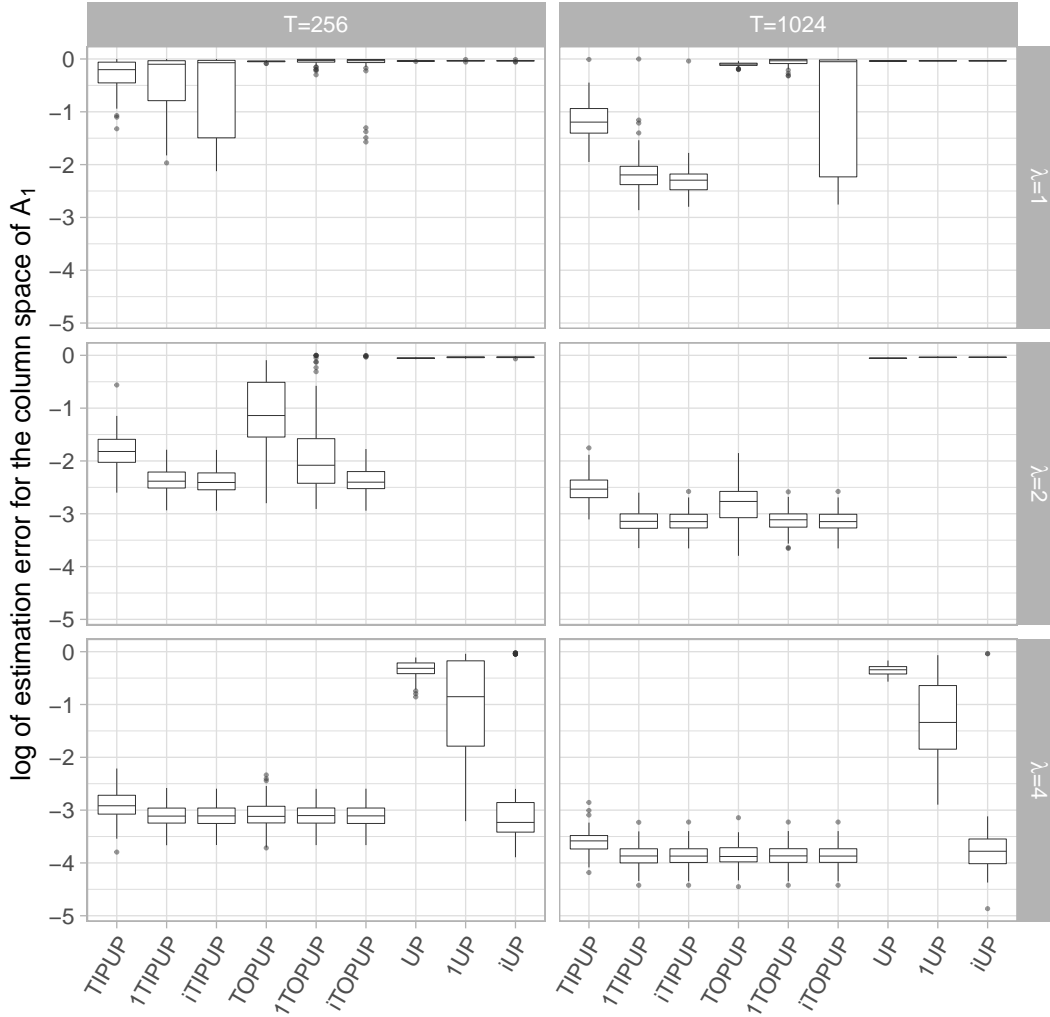


FIG 1. *Experiment 1 under Configuration I. Boxplot of the logarithm of the estimation error of A_1 . Nine methods (i)-(ix) are considered in total. Three rows correspond to three signal-to-noise strengths $\lambda = 1, 2, 4$. Two columns correspond to two sample sizes $T = 256, 1024$.*

repetitions are connected, for the $T = 256$ case. The top two panels show that when signal is weak and the sample size is small, the initial estimates may be poor, and the iterative methods may need certain accuracy in the initial estimates to produce further improvement. This reemphasizes the condition on the initial estimate in the theorems. The bottom two panels show that when signal is stronger, the relatively more accurate initial estimates enable the iterative methods to improve the estimates. Again, TOPUP initial estimates are not as accurate as the TIPUP estimates.

Experiment 2 under Configuration I. We use the same setting as above, but vary $h_0 = 1, 2, \dots, 5$, and fix $\hat{r}_2 = 1$ in the process of the estimation. Figure 3 provides the results. It can be seen that when there is no signal cancellation, the choice of h_0 does not affect the performance dramatically. When h_0 increases, the performance of all TIPUP-based and TOPUP-based algorithms becomes slightly worse most of the time.

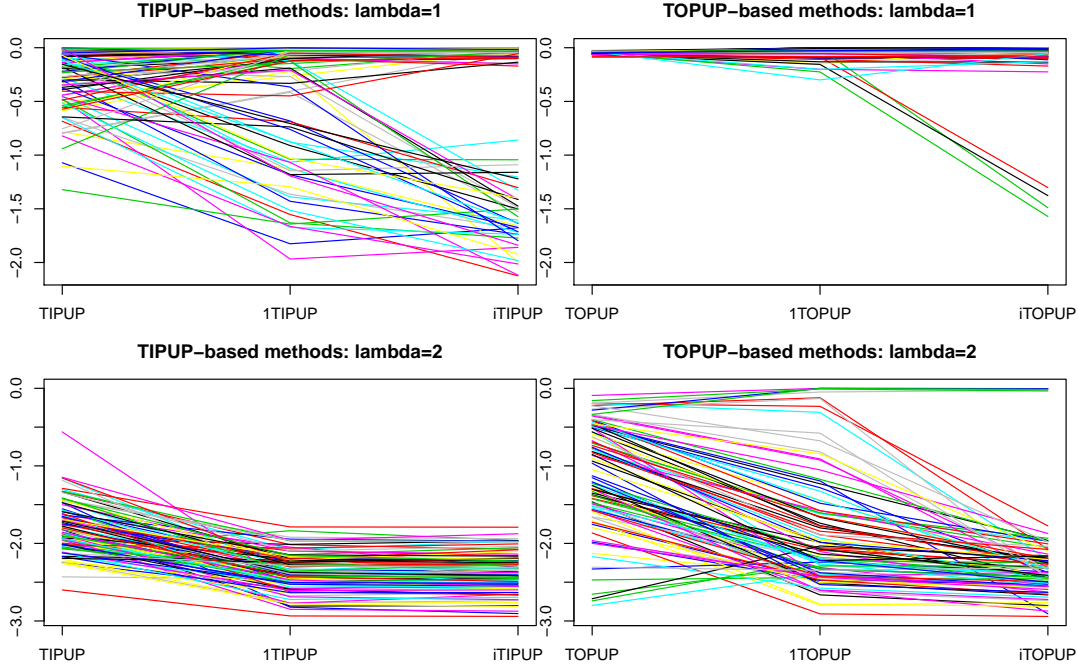


FIG 2. Experiment 1 under Configuration I. Trajectory of the logarithm of the estimation error of A_1 with fixed sample size $T = 256$. Two rows correspond to two signal-to-noise strengths $\lambda = 1, 2$. Two columns correspond to TIPUP-based and TOPUP-based methods respectively.

Experiment 3 under Configuration I. This experiment is conducted to verify the bounds on the sample size for iTOPUP in Corollary 3.3 and for iTIPUP in Corollary 3.4. We set the off-diagonal entries of the covariance matrices of the noise as $\psi_1 = \psi_2 = 0.4$ and the AR coefficient of the factor time series as $\phi = 0.9$, and vary the sample size $T = 16, 64, 256, 1024, 4096$ and signal strength $\lambda = 1, 2, 4, 8$. Again, we use $h_0 = 1$ and $\hat{r}_2 = 1$ in the process of the estimation.

Table 1 provides the values of δ_0, δ_1 in Assumption 3 as signal strength λ varies and the lower bounds on the sample size in (3.25) and (3.27) when the values of δ_0, δ_1 are plugged in. We have $\delta_2 = 0$ in this case and $d = d_1 \times d_2 = 256$. Figure 4 shows the simulation results, where the raw estimation error instead of the logarithm is given. When the raw errors are close to 1, it implies the estimation is not accurate. It is corroborated that the sample size T can be much smaller than d in the strong factor case.

λ	δ_0	δ_1	iTOPUP bound (3.25)	iTIPUP bound (3.27)
1	1.00	1.00	4096.00	256.00
2	1.00	1.00	4096.00	256.00
4	0.82	0.86	829.44	81.00
8	0.57	0.61	51.84	5.06

TABLE 1

Experiment 3 under Configuration I. This table provides the values of δ_0, δ_1 in Assumption 3 as signal strength λ varies and the lower bounds on the sample size in (3.25) and (3.27) when the values of δ_0, δ_1 are plugged in.

Under Configuration II, we performed two experiments: the two AR coefficients for the two independent univariate time series f_{1t} and f_{2t} are $\phi_1 = 0.8$ and $\phi_2 = 0.6$ in Experiment

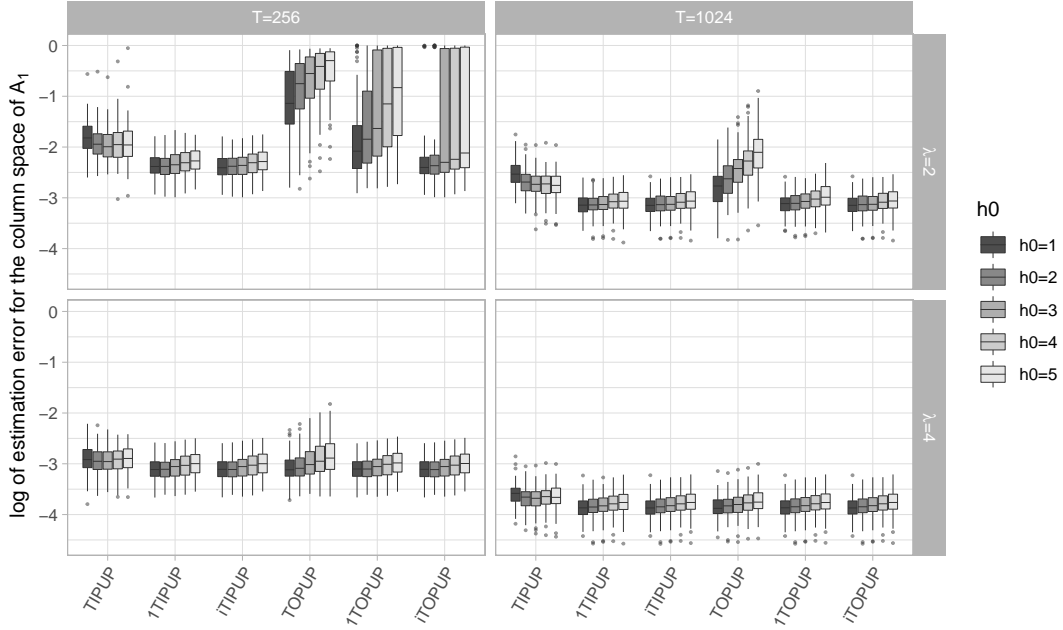


FIG 3. Experiment 2 under Configuration I. Boxplot of the logarithm of the estimation error of A_1 . Six methods (iv)-(ix) with five choices of h_0 are considered in total. Two rows correspond to two signal-to-noise strengths $\lambda = 2, 4$. Two columns correspond to two sample sizes $T = 256, 1024$.

1 and $\phi_1 = 0.8$ and $\phi_2 = -0.8$ in Experiment 2. Experiment 1 under configuration II satisfies Assumption 2(i)-(ii) because there is no signal cancellation; Experiment 2 under configuration II satisfies Assumption 2(i) for TOPUP related methods. When $h_0 = 1$, Experiment 2 does not satisfy Assumption 2(ii) for TIPUP related methods as there is a severe signal cancellation. However, using $h_0 = 2$ significantly reduces signal cancellation as lag 2 auto-cross-covariance does not cancel each other in TIPUP related methods.

Experiment 1 under Configuration II. Figure 5 shows the boxplot of the logarithm of the estimation errors of 8 methods including (iv)-(ix) and mixed (x)-(xi) with TIPUP initiation and TOPUP iteration. Again, the performance of the mixed (xii)-(xiii) procedures with iTIPUP iteration is not as good as that of iTIPUP hence not shown. Here we use different sample sizes, with the signal strength fixed at $\lambda = 1$ and two h_0 values: $h_0 = 1$ and $h_0 = 2$. The theoretical λ_1 defined in (3.5) and λ_1^* in (3.6) under the stationary auto-cross-moments of the factor process are given in the figure. Note that they are different for different h_0 . It shows that the mixed TIPUP-1TOPUP method can slightly improve 1TOPUP because of the better initialization. With larger sample size $T = 1024$, TIPUP-1TOPUP also slightly outperforms 1TIPUP. In this case, using the larger $h_0 = 2$ provides slightly poorer performance than $h_0 = 1$, as the lag-2 autocorrelation is significantly smaller than that of lag 1 for the underlying AR(1) process with $\phi_2 = 0.6$. The extra term adds limited signal, shown by the small differences in λ_1 and λ_1^* , but incorporates extra noise terms in the estimators. To see more clearly the impact of h_0 , we show the boxplots of the estimated λ_1^* and λ_1 using iTIPUP and iTOPUP, respectively, for $h_0 = 1, 2$ and 3, under different sample sizes in Figure 6. The theoretical values are marked with a diamond. It is seen that the estimated values are relatively close to the theoretical values. More importantly, they decrease as h_0 increases in this no-signal cancellation case.

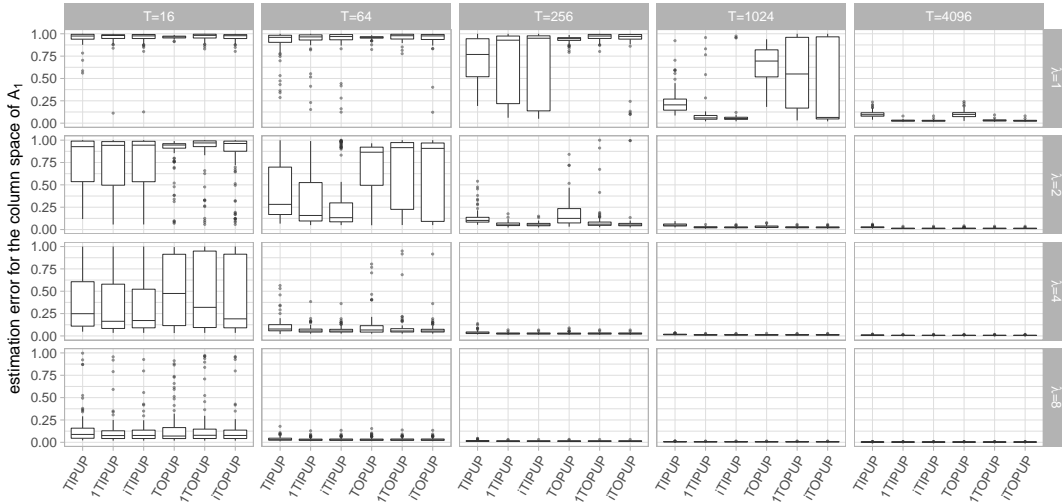


FIG 4. Experiment 3 under Configuration I. Boxplot of the estimation error of A_1 . Six methods (iv)-(ix) are considered in total. Four rows correspond to four signal-to-noise strengths $\lambda = 1, 2, 4, 8$. Five columns correspond to five sample sizes $T = 16, 64, 256, 1024, 4096$. This figure corroborates the theoretical lower bounds on the sample size in (3.25) and (3.27).

Experiment 2 under Configuration II. When $\phi_1 = 0.8$ and $\phi_2 = -0.8$, we can readily check that $\mathbb{E}(F_t F_{t-1}^\top) = (\phi_1 + \phi_2)\sigma^2 = 0$. Therefore, in the TIPUP-related procedure for estimating A_1 with $h_0 = 1$, the signal completely cancels out. Since the ranks r_1 and r_2 are fixed, we have $\delta_2 = \infty$ for $h_0 = 1$, and the corresponding $\lambda_1^* = 0$. Figure 7 shows the boxplot of the logarithm of the estimation error of A_1 for 8 methods including (iv)-(ix) and mixed (x)-(xi) with two choices of $h_0 = 1$ and $h_0 = 2$. We fix the signal strength to be $\lambda = 1$ to isolate the effect of h_0 . When $h_0 = 1$, both initialization TIPUP and TOPUP do not perform well. But 1TOPUP and iTOPUP improve the performance of TOPUP significantly with TOPUP iteration while 1TIPUP and iTIPUP cannot improve TIPUP. This is because signal cancellation has significant impact on TIPUP based procedures while having no impact on TOPUP based procedures. To our pleasant surprise, when $h_0 = 1$, the mixed TIPUP-1TOPUP is better than both 1TIPUP and 1TOPUP, and the mixed TIPUP-iTOPUP is similar to iTOPUP and much better than iTIPUP. When using $h_0 = 2$, the noise cancellation is mild and $(\lambda_1^*)^2 = 1.78$. Since r_k are fixed, we have $\delta_2 < \infty$. Note that in this case the signal using TIPUP only comes from lag-2 cross product and is weaker than that using TOPUP related procedures. The difference does not have impact on the convergence rate, but on the signal to noise ratio. Comparing the left two subfigures with the right ones of Figure 7, it is seen that using $h_0 = 2$ always boosts the performance of TIPUP-related methods significantly. Meanwhile, the TOPUP based methods are not sensitive to the choice of h_0 . When $h_0 = 2$, the non-iterative TIPUP performs better than TOPUP, 1TIPUP performs better than 1TOPUP, but after convergence, iTOPUP performs better than iTIPUP. Because the initialization TIPUP is better than TOPUP for $h_0 = 2$, it is of no surprise to see that TIPUP-1TOPUP behaves better than 1TIPUP and 1TOPUP, and TIPUP-iTOPUP is similar as iTOPUP and slightly better than iTIPUP.

Again, to see more clearly the impact of h_0 in this case with noise cancellation, we show the boxplots of the estimated λ_1^* and λ_1 using iTIPUP and iTOPUP, respectively, for $h_0 = 1, 2$ and 3 in Figure 8. It is seen that the iTOPUP procedure remains robust in estimating λ_1 under the noise-cancellation case. And λ_1 decreases as h_0 increases. However, iTIPUP is

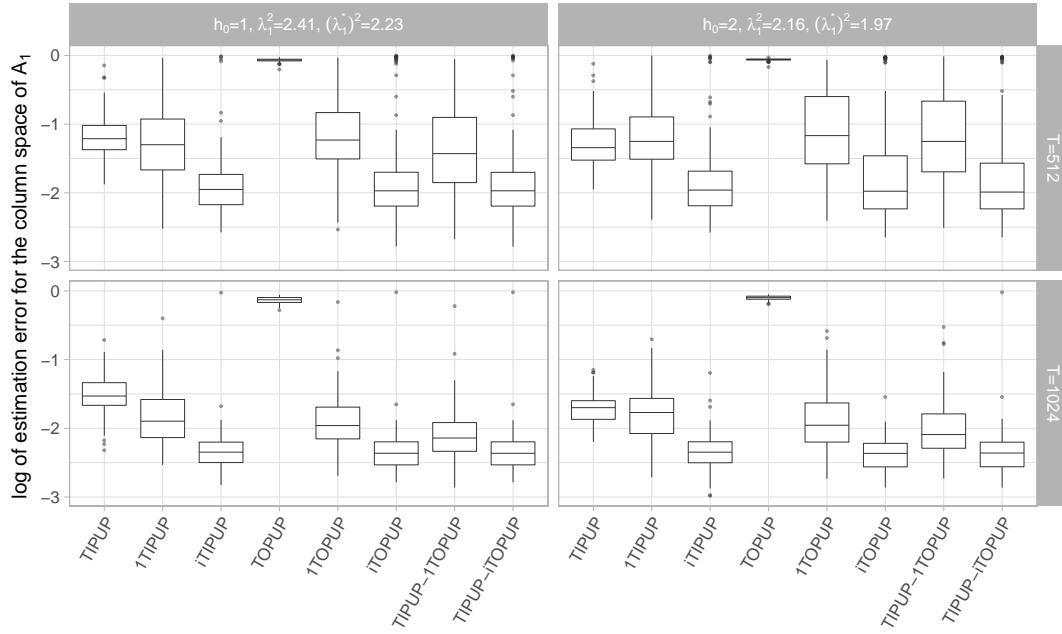


FIG 5. Experiment 1 under Configuration II. Boxplot of the logarithm of the estimation error of A_1 . Eight methods are considered in total. Two rows correspond to two sample sizes $T = 512, 1024$. Two columns correspond to two choices of h_0 . The population signal strengths λ_1^2 (3.5) and λ_1^{*2} (3.6) for different h_0 are provided on the top.

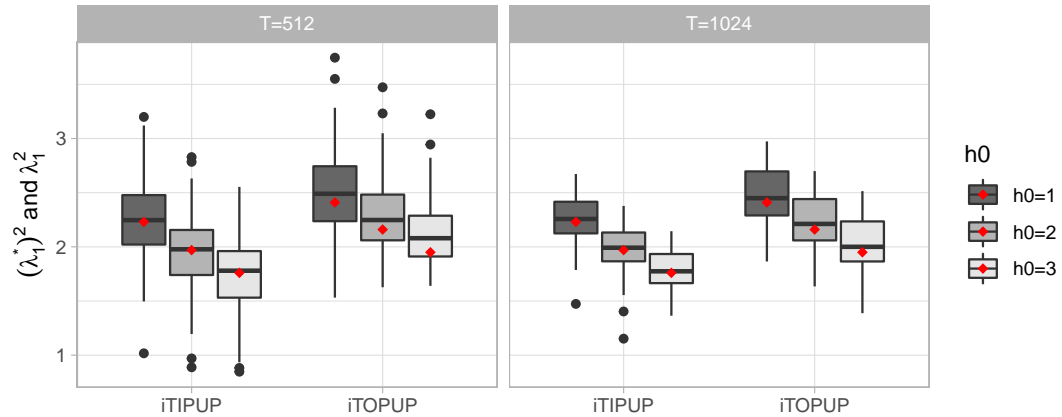


FIG 6. Experiment 1 under Configuration II. Boxplot of the sample estimates of the signal strengths λ_1^2 (3.5) and λ_1^{*2} (3.6) over 100 replications for iTIPUP and iTOPUP with three choices of h_0 . Two panels correspond to two sample sizes $T = 512, 1024$. The superimposed red diamonds are the population version of the signal strengths.

very different. Although when using $h_0 = 1$ the estimated λ_1^* significantly overestimates the theoretical value $\lambda_1^* = 0$, they are still much less than those from using $h_0 = 2$ and 3. The reversed order of the magnitude of λ_1^* as h_0 increases can be potentially used to detect signal cancellation in practice, though the theoretical property of the estimators of λ_1^* (e.g. standard deviation) is technically challenging to obtain. In practice, when one observes such

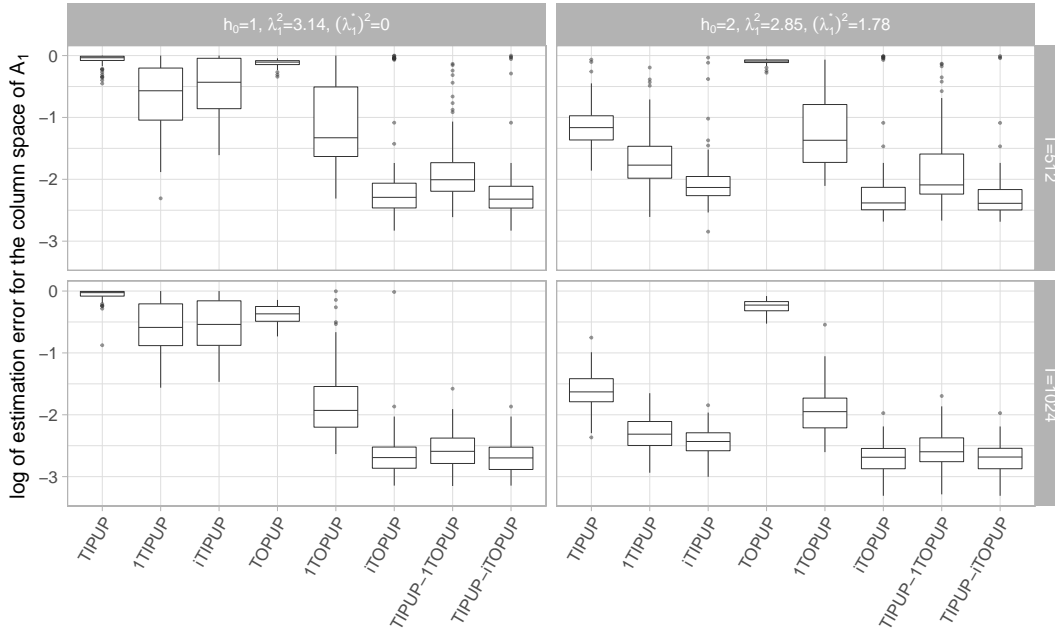


FIG 7. Experiment 2 under Configuration II. Boxplot of the logarithm of the estimation error of A_1 . 8 methods are considered in total. Two rows correspond to two sample sizes $T = 512, 1024$. Two columns correspond to two choices of h_0 . The population signal strengths λ_1^2 (3.5) and λ_1^{*2} (3.6) for different h_0 are provided on the top.

a reversed order, it is recommended to use iTOPUP as a conservative estimator. Of course, the behaviors of λ_k and λ_k^* depend on the auto-cross-moment structure of the underlying factor process. For example, if the factor process follows a MA(2) model with zero lag-1 autocorrelation ($f_t = e_t + \theta_2 e_{t-2}$), then λ_1 and λ_1^* under $h_0 = 2$ would be larger than those under both $h_0 = 1$ and $h_0 = 3$. But we expect that the pattern of λ_1 under different h_0 would be similar to that of λ_1^* under different h_0 , if there is no severe signal cancellation. Severe signal cancellation would make the patterns different.

Experiment under Configuration III. With order-3 tensor factor model (A.2), Configuration III satisfies Assumption 2, and Figure 9 shows the results. The message is almost the same as in Figure 3 for the matrix factor model. That is, TIPUP offers better initialization than TOPUP, which supports the theoretically smaller requirement on the sample size by TIPUP; iterative methods are better than one-step methods, which are in turn better than the non-iterative methods; when there is no signal cancellation, the iterative methods are in general not sensitive to the choice of h_0 and the non-iterative and one-step methods tend to behave slightly worse with larger values of h_0 ; when iTOPUP converges, its performance is better than that of iTIPUP as shown in the bottom right panel, which also verifies the theoretical claim of faster convergence rate of iTOPUP.

APPENDIX B: PROOF OF THEOREM 3.2

We focus on the case of $K = 2$ as the iTIPUP begins with mode- k matrix unfolding. In particular, we sometimes give explicit expressions only in the case of $k = 1$ and $K = 2$. For $K = 2$, we observe a matrix time series with $X_t = A_1 F_t A_2^T + E_t \in \mathbb{R}^{d_1 \times d_2}$. Recall that under the conditional expectation \mathbb{E} , F_1, \dots, F_T are fixed. Let U_1, U_2 be the left singular matrices of A_1 and A_2 respectively with $r_k = \text{rank}(U_k) = \text{rank}(A_k)$.

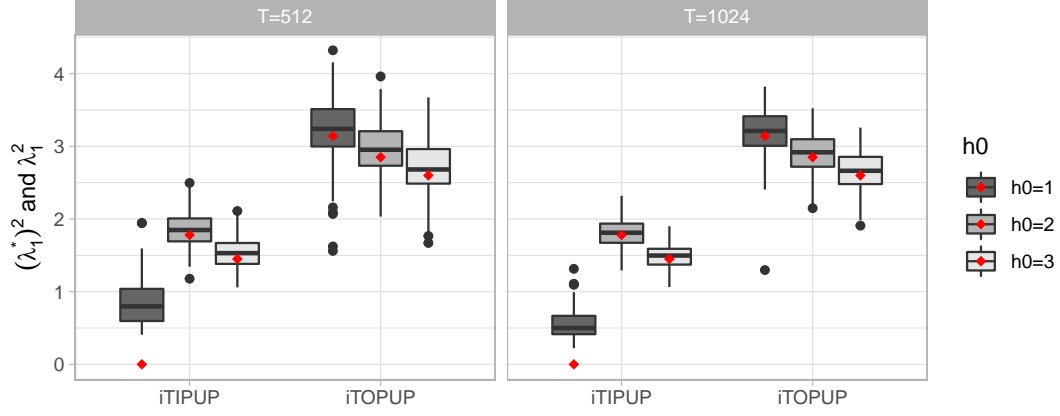


FIG 8. Experiment 2 under Configuration II. Boxplot of the sample estimates of the signal strengths λ_1^2 (3.5) and λ_1^{*2} (3.6) over 100 replications for iTIPUP and iTOPUP with three choices of h_0 . Two panels correspond to two sample sizes $T = 512, 1024$. The superimposed red diamonds are the population version of the signal strengths.

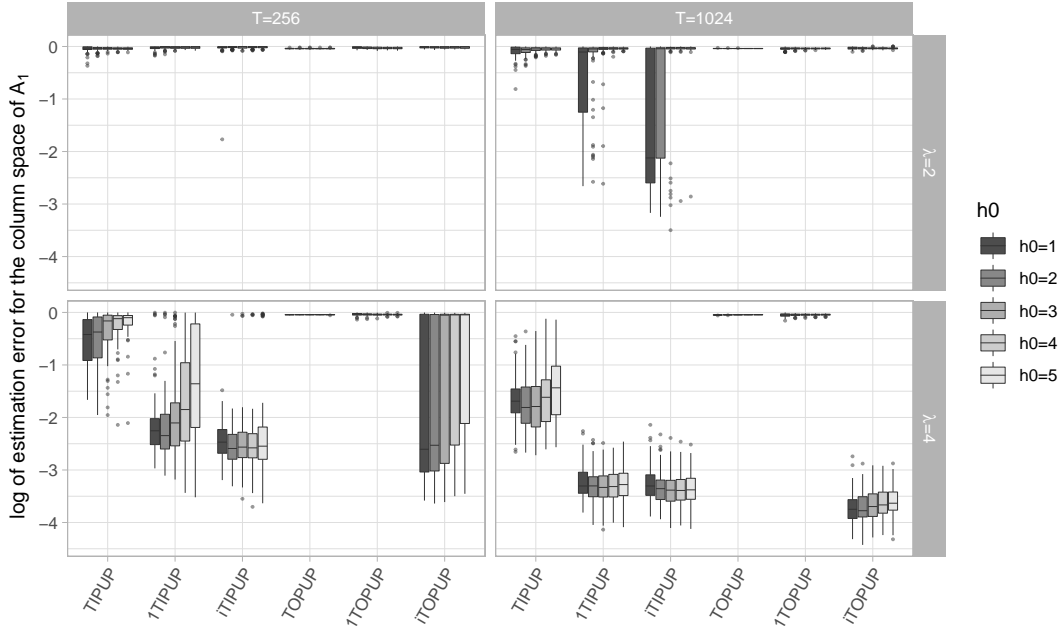


FIG 9. Experiment under Configuration III for order-3 tensor factor model. Boxplot of the logarithm of the estimation error of A_1 . Six methods (iv)-(ix) with five choices of h_0 are considered in total. Two rows correspond to two signal-to-noise strengths $\lambda = 2, 4$. Two columns correspond to two sample sizes $T = 256, 1024$.

We outline the proof as follows. Let $L_k^{(m)}$ be the loss (2.1) for $\hat{U}_k^{(m)}$ or equivalently the spectral norm error for $\hat{P}_k^{(m)} = \hat{U}_k^{(m)} \hat{U}_k^{(m)\top}$, $k = 1, \dots, K$, and $L^{(m)}$ their maximum,

$$(B.1) \quad L_k^{(m)} = \|\hat{P}_k^{(m)} - P_k\|_S, \quad L^{(m)} = \max_{k=1,2,\dots,K} L_k^{(m)}.$$

From [Chen, Yang and Zhang \(2019\)](#), $\overline{\mathbb{E}}[L_k^{(0)}] \lesssim R_k^{*(0)}$ as we mentioned in (3.14). By applying the Gaussian concentration inequality for Lipschitz functions and Lemme G.2 in their analysis, we have

$$(B.2) \quad L^{(0)} \leq C_1^{(\text{TIPUP})} R^{*(0)} \quad \text{with} \quad R^{*(0)} = \max_{1 \leq k \leq K} R_k^{*(0)}$$

in an event Ω_0 with $\overline{\mathbb{P}}(\Omega_0) \geq 1 - 5^{-1} \sum_{k=1}^K e^{-d_k}$. This is similar to (B.10) below.

After the initialization with $\hat{U}_k^{(0)}$, the algorithm iteratively produces estimates $\hat{U}_k^{(m)}$ from $m = 1$ to $m = J$. Define the matrix-valued operator $\text{TIPUP}_1(\cdot)$ as

$$\text{TIPUP}_1(\tilde{U}_2) = \left(\sum_{t=h+1}^T \frac{X_{t-h} \tilde{U}_2 \tilde{U}_2^\top X_t^\top}{T-h}, h = 1, \dots, h_0 \right) \in \mathbb{R}^{d_1 \times (d_1 h_0)}$$

for any matrix-valued variable $\tilde{U}_2 \in \mathbb{R}^{d_2 \times r_2}$. Given $\hat{U}_2^{(m)}$, the $(m+1)$ -th iteration produces estimates

$$\hat{U}_1^{(m+1)} = \text{LSVD}_{r_1}(\text{TIPUP}_1(\hat{U}_2^{(m)})), \quad \hat{P}_1^{(m+1)} = \hat{U}_1^{(m+1)} \hat{U}_1^{(m+1)\top}.$$

The ‘‘noiseless’’ version of this update is given by

$$(B.3) \quad \Theta_{1,h}^*(\tilde{U}_2) = \sum_{t=h+1}^T \frac{A_1 F_{t-h} A_2^\top \tilde{U}_2 \tilde{U}_2^\top A_2 F_t^\top A_1^\top}{T-h}, \quad \overline{\mathbb{E}}[\text{TIPUP}_1](\tilde{U}_2) = \Theta_{1,1:h_0}^*(\tilde{U}_2),$$

with $\Theta_{1,1:h_0}^*(\tilde{U}_2) = (\Theta_{1,h}^*(\tilde{U}_2), h = 1, \dots, h_0)$ as in (3.4), giving error free ‘‘estimates’’,

$$U_1 = \text{LSVD}_{r_1}(\overline{\mathbb{E}}[\text{TIPUP}_1](\hat{U}_2^{(m)})), \quad P_1 = U_1 U_1^\top,$$

when $\overline{\mathbb{E}}[\text{TIPUP}_1](\hat{U}_2^{(m)})$ is of rank r_1 . Thus, by Wedin’s theorem ([Wedin \(1972\)](#)),

$$L_1^{(m+1)} = \|\hat{P}_1^{(m+1)} - P_1\|_S \leq \frac{2\|\text{TIPUP}_1(\hat{U}_2^{(m)}) - \overline{\mathbb{E}}[\text{TIPUP}_1](\hat{U}_2^{(m)})\|_S}{\sigma_{r_1}(\overline{\mathbb{E}}[\text{TIPUP}_1](\hat{U}_2^{(m)})}).$$

We note that $\hat{U}_2^{(m)}$ is plugged-in after the conditional expectation in $\overline{\mathbb{E}}[\text{TIPUP}_1](\cdot)$. For general $1 \leq k \leq K$, we define $\text{TIPUP}_k(\tilde{U}_{-k})$ and $\overline{\mathbb{E}}[\text{TIPUP}_k](\tilde{U}_{-k})$ as matrix-valued functions of $\tilde{U}_{-k} = \odot_{j \neq k} \tilde{U}_j$. We will prove by induction that $\sigma_{r_k}(\overline{\mathbb{E}}[\text{TIPUP}_k](\hat{U}_{-k}^{(m)}))$, the denominator in the above inequality, is no smaller than a half of its ideal version as in (3.6), e.g.

$$(B.4) \quad 2\sigma_{r_1}(\overline{\mathbb{E}}[\text{TIPUP}_1](\hat{U}_2^{(m)})) \geq \sigma_{r_1}(\overline{\mathbb{E}}[\text{TIPUP}_1](U_2)) = h_0^{1/2} \lambda_1^{*2},$$

in the case of $k = 1$ and $K = 2$. It would then follow that

$$(B.5) \quad L_1^{(m+1)} = \|\hat{P}_1^{(m+1)} - P_1\|_S \leq \frac{\|\text{TIPUP}_1(\hat{U}_2^{(m)}) - \overline{\mathbb{E}}[\text{TIPUP}_1](\hat{U}_2^{(m)})\|_S}{h_0^{1/2} \lambda_1^{*2}/4}.$$

To bound the numerator on the right-hand side of (B.5), we write

$$(B.6) \quad \text{TIPUP}_1(\tilde{U}_2) - \overline{\mathbb{E}}[\text{TIPUP}_1](\tilde{U}_2) = \sum_{j=1}^3 \left(\Delta_{j,1,h}^*(\tilde{U}_2 \tilde{U}_2^\top), h = 1, \dots, h_0 \right) \in \mathbb{R}^{d_1 \times (d_1 h_0)}$$

as both $\text{TIPUP}_1(\tilde{U}_2)$ and $\overline{\mathbb{E}}[\text{TIPUP}_1](\tilde{U}_2)$ are linear in $\tilde{U}_2 \tilde{U}_2^\top$, where for any $\tilde{M}_2 \in \mathbb{R}^{d_2 \times d_2}$

$$\Delta_1^*(\tilde{M}_2) := \Delta_{1,1,h}^*(\tilde{M}_2) := \sum_{t=h+1}^T A_1 F_{t-h} A_2^\top \tilde{M}_2 E_t^\top / (T-h),$$

$$\Delta_2^*(\tilde{M}_2) := \Delta_{2,1,h}^*(\tilde{M}_2) := \sum_{t=h+1}^T E_{t-h} \tilde{M}_2 A_2 F_t^\top A_1^\top / (T-h),$$

$$\Delta_3^*(\tilde{M}_2) := \Delta_{3,1,h}^*(\tilde{M}_2) := \sum_{t=h+1}^T E_{t-h} \tilde{M}_2 E_t^\top / (T-h).$$

As $\Delta_{j,1,h}^*(\widetilde{M}_2)$ is linear in \widetilde{M}_2 , the numerator on the right-hand of (B.5) can be bounded by

$$(B.7) \quad \begin{aligned} & \|\text{TIPUP}_1(\widehat{U}_2^{(m)}) - \overline{\mathbb{E}}[\text{TIPUP}_1](\widehat{U}_2^{(m)})\|_S \\ & \leq \|\text{TIPUP}_1(U_2) - \overline{\mathbb{E}}[\text{TIPUP}_1](U_2)\|_S + L^{(m)}(2K-2) \sum_{j=1}^3 h_0^{1/2} \max_{h \leq h_0} \|\Delta_{j,1,h}^*\|_{1,S,S} \end{aligned}$$

with an application of Cauchy-Schwarz inequality for the sum over $h = 1, \dots, h_0$, where $\|\Delta_{j,1,h}^*\|_{1,S,S}$ are norms of the $\mathbb{R}^{d_2 \times d_2} \rightarrow \mathbb{R}^{d_1 \times d_1}$ linear mappings $\Delta_{j,1,h}^*$ defined as

$$\|\Delta_j^*\|_{1,S,S} = \|\Delta_{j,1,h}^*\|_{1,S,S} := \max_{\|\widetilde{M}_2\|_S \leq 1, \text{rank}(\widetilde{M}_2) \leq r_2} \|\Delta_{j,1,h}^*(\widetilde{M}_2)\|_S.$$

For general $1 \leq k \leq K$, $\Delta_{j,k,h}^*$ is an $\mathbb{R}^{d_{-k} \times d_{-k}} \rightarrow \mathbb{R}^{d_k \times d_k}$ mapping, and (G.2) of Lemma G.1 (iii) gives the general version of (B.7) with

$$\|\Delta_{j,k,h}^*\|_{k,S,S} := \max_{\|\widetilde{M}_\ell\|_S \leq 1, \text{rank}(\widetilde{M}_\ell) \leq r_\ell, \forall \ell \neq k} \|\Delta_{j,k,h}^*(\odot_{\ell \neq k} \widetilde{M}_\ell)\|_S$$

because it is applied to $\widetilde{M}_{-k} = \odot_{\ell \neq k} \widehat{P}_\ell^{(m)} - \odot_{\ell \neq k} P_\ell$ with $\|\widehat{P}_\ell^{(m)} - P_\ell\|_S \leq L^{(m)}$.

We claim that in certain events $\Omega_j, j = 1, 2, 3$, with $\overline{\mathbb{P}}(\Omega_j) \geq 1 - 5^{-1} \sum_{k=1}^K e^{-d_k}$,

$$(B.8) \quad \|\Delta_{j,k,h}^*\|_{k,S,S} \leq \rho \lambda_k^{*2} / (24(K-1)), \quad \forall 1 \leq k \leq K.$$

For simplicity, we will only prove this inequality for $k = 1$ and $K = 2$ in the form of

$$(B.9) \quad \overline{\mathbb{P}} \left\{ \|\Delta_j^*\|_{1,S,S} = \max_{\|\widetilde{M}_2\|_S \leq 1, \text{rank}(\widetilde{M}_2) \leq r_2} \|\Delta_{j,1,h}^*(\widetilde{M}_2)\|_S \geq \rho \lambda_1^{*2} / 24 \right\} \leq 5^{-1} e^{-d_2}$$

as the proof of its counter part for general $1 \leq k \leq K$ is similar.

Define the ideal version of the ratio in (B.5) for general $1 \leq k \leq K$ as

$$L_k^{(\text{ideal})} = \frac{\|\text{TIPUP}_k(U_{-k}) - \overline{\mathbb{E}}[\text{TIPUP}_k](U_{-k})\|_S}{h_0^{1/2} \lambda_k^{*2} / 4}, \quad L^{(\text{ideal})} = \max_{1 \leq k \leq K} L_k^{(\text{ideal})}.$$

As $U_{-k} = \odot_{j \neq k} U_j$ is true and deterministic, by (3.15) and (B.6), the proof of (B.8) also implies

$$(B.10) \quad L^{(\text{ideal})} \leq C_1^{(\text{TIPUP})} R^{*(\text{ideal})} \quad \text{with} \quad R^{*(\text{ideal})} = \max_{1 \leq k \leq K} R_k^{*(\text{ideal})}$$

in an event Ω_4 with $\overline{\mathbb{P}}(\Omega_4) \leq 5^{-1} \sum_{k=1}^K e^{-d_k}$, where $R_k^{*(\text{ideal})}$ is as in (3.15).

Putting together (B.1), (B.5), (B.7) and (B.8), we find that in the event $\cap_{j=0}^4 \Omega_j$

$$(B.11) \quad L_k^{(m+1)} \leq L_k^{(\text{ideal})} + L^{(m)} \frac{(6K-6) \max_{j,k,h} \|\Delta_{j,k,h}^*\|_{k,S,S}}{\lambda_k^{*2} / 4} \leq L_k^{(\text{ideal})} + \rho L^{(m)},$$

which implies by induction

$$(B.12) \quad L^{(m+1)} \leq (1 + \dots + \rho^m) L^{(\text{ideal})} + \rho^{m+1} L^{(0)},$$

and then the conclusions follows from (B.2) and (B.10). We divide the rest of the proof into 4 steps to prove (B.9) for $j = 1, 2, 3$ and (B.4).

Step 1. We prove (B.9) for the $\Delta_1^*(\widetilde{M}_2)$ in (B.6). By Lemma G.1 (ii), there exist $\widetilde{M}^{(\ell, \ell')} \in \mathbb{R}^{d_2 \times d_2}$ of the form $W_\ell Q_{\ell'}^\top$ with $W_\ell \in \mathbb{R}^{d_2 \times r_2}$, $Q_{\ell'} \in \mathbb{R}^{d_2 \times r_2}$, $1 \leq \ell, \ell' \leq N_{d_2 r_2, 1/8} := 17^{d_2 r_2}$, such that $\|\widetilde{M}^{(\ell, \ell')}\|_S \leq 1$, $\text{rank}(\widetilde{M}^{(\ell, \ell')}) \leq r_2$ and

$$\|\Delta_1^*\|_{1,S,S} = \|\Delta_{1,1,h}^*\|_{1,S,S} \leq 2 \max_{\ell, \ell' \leq N_{d_2 r_2, 1/8}} \left\| \sum_{t=h+1}^T \frac{A_1 F_{t-h} A_2^\top \widetilde{M}^{(\ell, \ell')} E_t^\top}{T-h} \right\|_S.$$

To bound $\|\Delta_{1,k,h}^*\|_{1,S,S}$ for general $k \leq K$, we just need to replace $\widetilde{M}^{(\ell, \ell')}$ by $\odot_{j \neq k} \widetilde{M}_j^{(\ell, \ell')}$ and $N_{d_2 r_2, 1/8}$ by $N_{d_{-k}^*, 1/(8K-8)}$ with $d_{-k}^* = \sum_{j \neq k} d_j r_j$ as in Lemma G.1 (iii). We apply the Gaussian concentration inequality to the right-hand side above. Elementary calculation shows that

$$\begin{aligned} & \left\| \left\| \sum_{t=h+1}^T A_1 F_{t-h} A_2^\top \widetilde{M}^{(\ell, \ell')} E_t^\top \right\|_S - \left\| \sum_{t=h+1}^T A_1 F_{t-h} A_2^\top \widetilde{M}^{(\ell, \ell')} E_t^{*\top} \right\|_S \right\| \\ & \leq \left\| (A_1 F_1 A_2^\top, \dots, A_1 F_{T-h} A_2^\top) \begin{pmatrix} \widetilde{M}^{(\ell, \ell')} (E_{h+1}^\top - E_{h+1}^{*\top}) \\ \vdots \\ \widetilde{M}^{(\ell, \ell')} (E_T^\top - E_T^{*\top}) \end{pmatrix} \right\|_S \\ & \leq \|(A_1 F_1 A_2^\top, \dots, A_1 F_{T-h} A_2^\top)\|_S^{1/2} \left\| \text{diag}(\widetilde{M}^{(\ell, \ell')}) \begin{pmatrix} E_{h+1}^\top - E_{h+1}^{*\top} \\ \vdots \\ E_T^\top - E_T^{*\top} \end{pmatrix} \right\|_S \\ & \leq \sqrt{T} \|\Theta_{1,0}^*\|_S^{1/2} \left\| \begin{pmatrix} E_{h+1}^\top - E_{h+1}^{*\top} \\ \vdots \\ E_T^\top - E_T^{*\top} \end{pmatrix} \right\|_F. \end{aligned}$$

That is, $\left\| \sum_{t=h+1}^T A_1 F_{t-h} A_2^\top \widetilde{M}^{(\ell, \ell')} E_t^\top \right\|_S$ is a $\sqrt{T} \|\Theta_{1,0}^*\|_S^{1/2}$ Lipschitz function of (E_1, \dots, E_T) . Employing similar arguments in the proof of Theorem 2 in Chen, Yang and Zhang (2019), we have

$$\mathbb{E} \left\| \sum_{t=h+1}^T \frac{A_1 F_{t-h} A_2^\top \widetilde{M}^{(\ell, \ell')} E_t^\top}{T-h} \right\|_S \leq \frac{\sigma(8Td_1)^{1/2}}{T-h} \|\Theta_{1,0}^*\|_S^{1/2}.$$

Then, by Gaussian concentration inequalities for Lipschitz functions,

$$\mathbb{P} \left(\left\| \sum_{t=h+1}^T \frac{A_1 F_{t-h} A_2^\top \widetilde{M}^{(\ell, \ell')} E_t^\top}{T-h} \right\|_S - \frac{\sigma(8Td_1)^{1/2}}{T-h} \|\Theta_{1,0}^*\|_S^{1/2} \geq \frac{\sigma\sqrt{T}}{T-h} \|\Theta_{1,0}^*\|_S^{1/2} x \right) \leq e^{-\frac{x^2}{2}}.$$

Hence,

$$\mathbb{P} \left(\|\Delta_1^*\|_{1,S,S}/2 \geq \frac{\sigma(8Td_1)^{1/2}}{T-h} \|\Theta_{1,0}^*\|_S^{1/2} + \frac{\sigma\sqrt{T}}{T-h} \|\Theta_{1,0}^*\|_S^{1/2} x \right) \leq N_{d_2 r_2, 1/8}^2 e^{-\frac{x^2}{2}}.$$

As $T \geq 4h_0$ and $K = 2$, this implies with $x \asymp \sqrt{d_2 r_2}$ that in an event with at least probability $1 - e^{-d_2/5}$,

$$\|\Delta_1^*\|_{1,S,S} \leq \frac{C_{1,K}^{(\text{iter})} \sigma T^{-1/2} \|\Theta_{1,0}^*\|_S^{1/2} (\sqrt{d_1} + \sqrt{d_2 r_2})}{24(K-1)} \leq \frac{C_{1,K}^{(\text{iter})} (1 + \sqrt{d_2 r_2/d_1}) \lambda_1^{*2} R_1^{*(\text{ideal})}}{24},$$

with the $R_k^{*(\text{ideal})}$ in (3.15) and a constant $C_{1,K}^{(\text{iter})}$ depending on K only. In this event, (3.17) gives $24(\lambda_k^*)^{-2} \|\Delta_{1,k,h}^*\|_{k,S,S} \leq C_{1,K}^{(\text{iter})} (R_k^{*(\text{ideal})} + R_k^{*(\text{add})}) \leq \rho$. Thus, (B.9) holds for $\Delta_1^*(\widetilde{M}_2)$.

Step 2. Inequality (B.9) for $\Delta_2^*(\widetilde{M}_2)$ follow from the same argument as the above step.

Step 3. Here we prove (B.9) for the $\Delta_3^*(\widetilde{M}_2)$ in (B.6). By Lemma G.1 (ii), we can find $U_2^{(\ell)}, U_2^{(\ell')} \in \mathbb{R}^{d_2 \times r_2}$, $1 \leq \ell, \ell' \leq N_{d_2 r_2, 1/8}$ such that $\|U_2^{(\ell)}\|_S \leq 1$, $\|U_2^{(\ell')}\|_S \leq 1$ and

$$(B.13) \quad \|\Delta_3^*\|_{1,S,S} = \|\Delta_{3,1,h}^*\|_{1,S,S} \leq 2 \max_{1 \leq \ell, \ell' \leq N_{d_2 r_2, 1/8}} \left\| \sum_{t=h+1}^T \frac{E_{t-h} U_2^{(\ell)} U_2^{(\ell')\top} E_t^\top}{T-h} \right\|_S.$$

We split the sum into two terms over the index sets, $S_1 = \{(h, 2h] \cup (3h, 4h] \cup \dots\} \cap (h, T]$ and its complement S_2 in $(h, T]$, so that $\{E_{t-h}, t \in S_a\}$ is independent of $\{E_t, t \in S_a\}$ for each $a = 1, 2$. Let $n_a = |S_a| r_2$. Define $G_a = (E_{t-h} U_2^{(\ell)}, t \in S_a) \in \mathbb{R}^{d_1 \times n_a}$ and $H_a = (E_t U_2^{(\ell')}, t \in S_a) \in \mathbb{R}^{d_1 \times n_a}$. Then, G_a, H_a are two independent Gaussian matrices. Note that

$$(B.14) \quad \left\| \sum_{t \in S_a} \frac{E_{t-h} U_2^{(\ell)} U_2^{(\ell')\top} E_t^\top}{T-h} \right\|_S = \left\| \frac{G_a H_a^\top}{T-h} \right\|_S.$$

Moreover, by Assumption 1, $\text{Var}(u^\top \text{vec}(G_a)) \leq \sigma^2$ and $\text{Var}(u^\top \text{vec}(H_a)) \leq \sigma^2$ for all unit vectors $u \in \mathbb{R}^{d_1 n_a}$, so that by Lemme G.2 (i),

$$\mathbb{P} \left\{ \|G_a H_a^\top\|_S / \sigma^2 \geq d_1 + 2\sqrt{d_1 n_a} + x(x + 2\sqrt{n_a} + 2\sqrt{d_1}) \right\} \leq e^{-x^2/2}, \quad x > 0.$$

As $\sum_{a=1}^2 n_a = r_2(T-h)$, it follows from (B.13), (B.14) and the above inequality that

$$\mathbb{P} \left\{ \frac{\|\Delta_{3,1,h}^*\|_{1,S,S}}{4\sigma^2} \geq \frac{(\sqrt{d_1} + x)^2}{T-h} + \frac{\sqrt{2r_2}(\sqrt{d_1} + x)}{\sqrt{T-h}} \right\} \leq 2N_{d_2 r_2, 1/8}^2 e^{-x^2/2}.$$

Thus, with $h_0 \leq T/4$, $x = \sqrt{d_1} + \sqrt{d_{-1}^*}$ and some constant $C_{1,K}^{(\text{iter})}$ depending on K only,

$$(B.15) \quad \|\Delta_{3,1,h}^*\|_{1,S,S} \leq \frac{(C_{1,K}^{(\text{iter})} - 1)(1 - h_0/T)^2 \sigma^2}{24(K-1)} \left(\frac{(\sqrt{d_1} + \sqrt{d_{-1}^*})\sqrt{r_{-1}}}{(T-h_0)^{1/2}} + \frac{(\sqrt{d_1} + \sqrt{d_{-1}^*})^2}{T-h_0} \right)$$

with at least probability $1 - e^{-d_1}/5$. As $\lambda_k^* \leq \|\Theta_{1,0}^*\|_S^{1/2} / (1 - h_0/T)^{1/2}$ by (3.6) and (3.3),

$$R_1^{*(\text{ideal})} \geq (\lambda_1^*)^{-1} (1 - h_0/T) \sigma (T - h_0)^{-1/2} \sqrt{d_1} + (\lambda_1^*)^{-2} \sigma T^{-1/2} \sigma \sqrt{d_1 r_{-1}}$$

by (3.15). Thus, in the event (B.15) and for $k = 1$ and $K = 2$,

$$\begin{aligned} & \|\Delta_{3,1,h}^*\|_{1,S,S} \\ & \leq \frac{C_{1,K}^{(\text{iter})} - 1}{24} \left(\left(1 + \sqrt{d_{-1}^*/d_1}\right) \lambda_1^{*2} R_1^{*(\text{ideal})} + \left(1 + \sqrt{d_{-1}^*/d_1}\right)^2 \lambda_1^{*2} (R_1^{*(\text{ideal})})^2 \right) \\ & \leq (\lambda_1^{*2}/24) (C_{1,K}^{(\text{iter})} - 1) (1 + R_k^{*(\text{ideal})} + R_k^{*(\text{add})}) (R_k^{*(\text{ideal})} + R_k^{*(\text{add})}) \end{aligned}$$

which is no greater than $\lambda_1^{*2} \rho / 24$ by the condition (3.17). This yields (B.9) for $\Delta_3^*(\widetilde{M}_2)$.

Step 4. Next, we prove (B.4) in the event $\cap_{j=0}^4 \Omega_j$. Note that,

$$\begin{aligned}
 \|\Theta_{1,h}^*(\widehat{U}_2^{(m)}) - \Theta_{1,h}^*(U_2)\|_S &= \left\| \sum_{t=h+1}^T \frac{A_1 F_{t-h} A_2^\top (\widehat{U}_2^{(m)} \widehat{U}_2^{(m)\top} - U_2 U_2^\top) A_2 F_t^\top A_1^\top}{T-h} \right\|_S \\
 &= \frac{1}{T-h} \left\| (A_1 F_1 A_2^\top, \dots, A_1 F_{T-h} A_2^\top) \begin{pmatrix} (\widehat{U}_2^{(m)} \widehat{U}_2^{(m)\top} - U_2 U_2^\top) A_2 F_1^\top A_1^\top \\ \vdots \\ (\widehat{U}_2^{(m)} \widehat{U}_2^{(m)\top} - U_2 U_2^\top) A_2 F_{T-h}^\top A_1^\top \end{pmatrix} \right\|_S \\
 &\leq \frac{1}{T-h} \left\| (A_1 F_1 A_2^\top, \dots, A_1 F_{T-h} A_2^\top) \right\|_S \left\| \text{diag} \left(\widehat{U}_2^{(m)} \widehat{U}_2^{(m)\top} - U_2 U_2^\top \right) \begin{pmatrix} A_2 F_1^\top A_1^\top \\ \vdots \\ A_2 F_{T-h}^\top A_1^\top \end{pmatrix} \right\|_S \\
 &\leq \|\widehat{U}_2^{(m)} \widehat{U}_2^{(m)\top} - U_2 U_2^\top\|_S \|\Theta_{1,0}^*\|_S / (1-h/T) \\
 &\leq L^{(m)} \|\Theta_{1,0}^*\|_S / (1-h_0/T).
 \end{aligned}$$

Hence, by the Cauchy-Schwarz inequality and (B.3),

$$\|\mathbb{E}[\text{TIPUP}_1](\widehat{U}_2^{(m)}) - \mathbb{E}[\text{TIPUP}_1](U_2)\|_S \leq \sqrt{h_0} L^{(m)} \|\Theta_{1,0}^*\|_S / (1-h_0/T).$$

By (3.6), $\lambda_1^{*2} h_0^{1/2} = \sigma_{r_1}(\text{mat}_1(\Theta_{1,1:h_0}^*)) = \sigma_{r_1}(\mathbb{E}[\text{TIPUP}_1](U_2))$. Thus, by Weyl's inequality,

$$\sigma_{r_1}(\mathbb{E}[\text{TIPUP}_1](\widehat{U}_2^{(m)})) \geq \lambda_1^{*2} h_0^{1/2} - 2\sqrt{h_0} L^{(m)} \|\Theta_{1,0}^*\|_S \geq \sigma_{r_1}(\mathbb{E}[\text{TIPUP}_1](U_2)) / 2 = \lambda_1^{*2} h_0^{1/2} / 2.$$

when $\min_{k \leq K} \lambda_k^{*2} \geq 4L^{(m)} \|\Theta_{1,0}^*\|_S$. We prove this condition by induction in the event $\cap_{j=0}^4 \Omega_j$. By (3.17) and (B.2), $4L^{(0)} \|\Theta_{1,0}^*\|_S \leq 4C_1^{(\text{TIPUP})} R^{*(0)} \|\Theta_{1,0}^*\|_S \leq \min_{k \leq K} \lambda_k^{*2}$. Given the induction assumption $4L^{(m)} \|\Theta_{1,0}^*\|_S \leq \min_{k \leq K} \lambda_k^{*2}$, (B.4) holds for the same m , so that (B.12), (B.10) and (B.2),

$$L^{(m+1)} \leq C_1^{(\text{TIPUP})} \{2(1 + \dots + \rho^m) R^{*(\text{ideal})} + \rho^{m+1} R^{*(0)}\} \leq C_1^{(\text{TIPUP})} 2(1 - \rho)^{-1} R^{*(0)}.$$

It then follows from (3.17) that $4L^{(m+1)} \|\Theta_{1,0}^*\|_S \leq C_1^{(\text{TIPUP})} 8(1 - \rho)^{-1} R^{*(0)} \|\Theta_{1,0}^*\|_S \leq \min_{k \leq K} \lambda_k^{*2}$. This completes the induction and the proof of the entire theorem.

APPENDIX C: PROOF OF THEOREM 3.1

Again, we focus on the case of $K = 2$ as the iTOPUP also begins with mode- k matrix unfolding. In particular, we sometimes give explicit expressions only in the case of $k = 1$ and $K = 2$. For $K = 2$, we observe a matrix time series with $X_t = A_1 F_t A_2^\top + E_t \in \mathbb{R}^{d_1 \times d_2}$. Recall $\mathbb{E}(\cdot) = \mathbb{E}(\cdot | \{\mathcal{F}_1, \dots, \mathcal{F}_T\})$. Let U_1, U_2 be the left singular matrices of A_1 and A_2 respectively with $r_k = \text{rank}(U_k) = \text{rank}(A_k)$. Recall \odot is kronecker product and \otimes is tensor product.

We outline the proof as follows, which has exactly the same structure as the proofs of Theorem 3.2.

Let $L_k^{(m)} = \|\widehat{P}_k^{(m)} - P_k\|_S$ and $L^{(m)} = \max_{k \leq K} L_k^{(m)}$ as in (B.1). From Chen, Yang and Zhang (2019), $\mathbb{E}[L_k^{(0)}] \lesssim R_k^{(0)}$ as we mentioned in (3.7). By applying the Gaussian concentration inequality for Lipschitz functions and Lemme G.2 in their analysis, we have

$$(C.1) \quad L^{(0)} \leq C_1^{(\text{TOPUP})} R^{(0)} \quad \text{with} \quad R^{(0)} = \max_{1 \leq k \leq K} R_k^{(0)}$$

in an event Ω_0 with $\mathbb{P}(\Omega_0) \geq 1 - 5^{-1} \sum_{k=1}^K e^{-d_k}$. This is similar to (C.8) below.

For any matrices $\tilde{U}_j \in \mathbb{R}^{d_j \times r_j}$ with $r_{K+j} = r_j$, define

$$(C.2) \quad \begin{aligned} & \text{TOPUP}_k(\tilde{U}_j, 1 \leq j \leq 2K, j \neq k, j \neq K+K) \\ &= \left(\text{mat}_k \left(\hat{\Sigma}_h \times_{j=1}^{k-1} \tilde{U}_j \times_{j=k+1}^{K+k-1} \tilde{U}_j \times_{j=K+k+1}^{2K} \tilde{U}_j \right) \right)_{h=1, \dots, h_0} \end{aligned}$$

and its noiseless version $\mathbb{E}[\text{TOPUP}_k]$. Here both TOPUP_k and $\mathbb{E}[\text{TOPUP}_k]$ are viewed as $\mathbb{R}^{d_k \times (d_k r_k^2 h_0)}$ -valued $(2K-2)$ -linear mappings with input \tilde{U}_j . When $\tilde{U}_{K+j} = \tilde{U}_j$ for all $j \neq k$, we write (C.2) as $\text{TOPUP}_k(\tilde{U}_{-k})$ and its noiseless version $\mathbb{E}[\text{TOPUP}_k](\tilde{U}_{-k})$. After the initialization with $\hat{U}_k^{(0)}$, the algorithm iteratively produces estimates $\hat{U}_k^{(m+1)}$ as the rank- r_k left singular matrix of $\text{TOPUP}_k(\hat{U}_{1:(k-1)}^{(m+1)}, \hat{U}_{(k+1):K}^{(m)})$. For $K=2$ and $k=1$,

$$\text{TOPUP}_1(\tilde{U}_2, \tilde{U}_4) = \text{mat}_1 \left(\sum_{t=h+1}^T \frac{X_{t-h} \tilde{U}_2 \otimes X_t \tilde{U}_4}{T-h}, h=1, \dots, h_0 \right) \in \mathbb{R}^{d_1 \times (d_1 r_2^2 h_0)}.$$

for any $\tilde{U}_j \in \mathbb{R}^{d_2 \times r_2}$, $j=2, 4$. Given $\hat{U}_2^{(m)}$, the $(m+1)$ -th iteration produces estimates

$$\hat{U}_1^{(m+1)} = \text{LSVD}_{r_1}(\text{TOPUP}_1(\hat{U}_2^{(m)})), \quad \hat{P}_1^{(m+1)} = \hat{U}_1^{(m+1)} \hat{U}_1^{(m+1)\top}.$$

When $\text{rank}(\mathbb{E}[\text{TOPUP}_1](\hat{U}_2^{(m)})) = r_1$, the “noiseless” version of this update is error free,

$$U_1 = \text{LSVD}_{r_1}(\mathbb{E}[\text{TOPUP}_1](\hat{U}_2^{(m)})), \quad P_1 = U_1 U_1^\top.$$

Thus, by Wedin’s theorem, (Wedin (1972)),

$$L_1^{(m+1)} = \|\hat{P}_1^{(m+1)} - P_1\|_S \leq \frac{2\|\text{TOPUP}_1(\hat{U}_2^{(m)}) - \mathbb{E}[\text{TOPUP}_1](\hat{U}_2^{(m)})\|_S}{\sigma_{r_1}(\mathbb{E}[\text{TOPUP}_1](\hat{U}_2^{(m)}))}.$$

We will prove by induction that $\sigma_{r_k}(\mathbb{E}[\text{TOPUP}_k](\hat{U}_{1:(k-1)}^{(m)}, \hat{U}_{(k+1):K}^{(m-1)}))$, the general version of the denominator on the right-hand side above, is no smaller than a half of its ideal version as in (3.5), e.g.

$$(C.3) \quad 2\sigma_{r_1}(\mathbb{E}[\text{TOPUP}_1](\hat{U}_2^{(m)})) \geq \sigma_{r_1}(\mathbb{E}[\text{TOPUP}_1](U_2)) = h_0^{1/2} \lambda_1^2,$$

in the case of $k=1$ and $K=2$. It would then follow that

$$(C.4) \quad L_1^{(m+1)} = \|\hat{P}_1^{(m+1)} - P_1\|_S \leq \frac{\|\text{TOPUP}_1(\hat{U}_2^{(m)}) - \mathbb{E}[\text{TOPUP}_1](\hat{U}_2^{(m)})\|_S}{h_0^{1/2} \lambda_1^2 / 4}.$$

To bound the numerator on the right-hand side of (C.4), we write

$$\Delta(\tilde{U}_2, \tilde{U}_4) = \text{TOPUP}_1(\tilde{U}_2, \tilde{U}_4) - \mathbb{E}[\text{TOPUP}_1](\tilde{U}_2, \tilde{U}_4)$$

and notice that for any orthonormal matrices $Q_j \in \mathbb{R}^{d_2 \times r_2}$

$$\|\Delta(Q_2 Q_2^\top \tilde{U}_2, Q_4 Q_4^\top \tilde{U}_4)\|_S \leq \|\Delta(Q_2, Q_4)\|_S \|Q_2^\top \tilde{U}_2\|_S \|Q_4^\top \tilde{U}_4\|_S$$

as the outer product is taken in (C.2). Moreover, because $\Delta(\tilde{U}_2, \tilde{U}_4)$ is bilinear,

$$\begin{aligned} \|\Delta(\tilde{U}_2, \tilde{U}_2)\|_S &\leq \|\Delta(P_2 \tilde{U}_2, P_2 \tilde{U}_2)\|_S + \|\Delta(P_2 \tilde{U}_2, P_2^\perp \tilde{U}_2)\|_S + \|\Delta(P_2^\perp \tilde{U}_2, \tilde{U}_2)\|_S \\ &\leq \|\Delta(U_2, U_2)\|_S + \|\Delta(U_2, Q_2)\|_S \|P_2^\perp \tilde{U}_2\|_S + \|\Delta(Q_2, \tilde{U}_2)\|_S \|P_2^\perp \tilde{U}_2\|_S, \end{aligned}$$

where $Q_2 = \text{LSVD}_{r_k}(P_2^\perp \tilde{U}_2)$ and $P_2^\perp = I_{d_2} - P_2$. Thus, due to $\|P_2^\perp \hat{U}_2^{(m)}\|_S \leq \|P_2 - \hat{P}_2^{(m)}\|_S \leq L^{(m)}$, the numerator on the right-hand of (C.4) can be bounded by

$$(C.5) \quad \|\text{TOPUP}_1(\hat{U}_2^{(m)}) - \mathbb{E}[\text{TOPUP}_1](\hat{U}_2^{(m)})\|_S \\ \leq \|\text{TOPUP}_1(U_2) - \mathbb{E}[\text{TOPUP}_1](U_2)\|_S + L^{(m)}(2K-2) \sum_{j=1}^3 h_0^{1/2} \max_{h \leq h_0} \|\Delta_{j,1,h}\|_{1,S,S},$$

where $\|\Delta\|_{1,S,S} := \max_{\|\tilde{U}_j\|_S \leq 1, \text{rank}(\tilde{U}_j) = r_2, j=2,4} \|\Delta(\tilde{U}_2, \tilde{U}_4)\|_S$ for any bilinear Δ and

$$\Delta_1(\tilde{U}_2, \tilde{U}_4) := \Delta_{1,1,h}(\tilde{U}_2, \tilde{U}_4) := \frac{1}{T-h} \sum_{t=h+1}^T \text{mat}_1(A_1 F_{t-h} A_2^\top \tilde{U}_2 \otimes E_t \tilde{U}_4), \\ \Delta_2(\tilde{U}_2, \tilde{U}_4) := \Delta_{2,1,h}(\tilde{U}_2, \tilde{U}_4) := \frac{1}{T-h} \sum_{t=h+1}^T \text{mat}_1(E_{t-h} \tilde{U}_2 \otimes A_1 F_t^\top A_2^\top \tilde{U}_4), \\ \Delta_3(\tilde{U}_2, \tilde{U}_4) := \Delta_{3,1,h}(\tilde{U}_2, \tilde{U}_4) := \frac{1}{T-h} \sum_{t=h+1}^T \text{mat}_1(E_{t-h} \tilde{U}_2 \otimes E_t \tilde{U}_4).$$

We claim that in certain events $\Omega_j, j = 1, 2, 3$, with $\bar{\mathbb{P}}(\Omega_j) \geq 1 - 5^{-1} \sum_{k=1}^K e^{-d_k}$, $\rho < 1$,

$$(C.6) \quad \|\Delta_{j,k,h}\|_{k,S,S} \leq \rho \lambda_k^2 / (24(K-1)), \quad \forall 1 \leq k \leq K,$$

in (C.5). For simplicity, we will only prove this inequality for $k = 1$ and $K = 2$,

$$(C.7) \quad \bar{\mathbb{P}} \left\{ \|\Delta_j\|_{1,S,S} = \max_{\|\tilde{M}_2\|_S \leq 1} \|\Delta_{j,1,h}(\tilde{U}_2, \tilde{U}_4)\|_S \geq \rho \lambda_1^2 / 24 \right\} \leq 5^{-1} e^{-d_2}.$$

Define the ideal version of the ratio in (C.4) for general $1 \leq k \leq K$ as

$$L_k^{(\text{ideal})} = \frac{4 \|\text{TOPUP}_k(U_{-k}) - \mathbb{E}[\text{TOPUP}_k](U_{-k})\|_S}{\sigma_{r_k}(\mathbb{E}[\text{TOPUP}_k](U_{-k}))}, \\ L^{(\text{ideal})} = \max_{1 \leq k \leq K} L_k^{(\text{ideal})}.$$

Note that $\sigma_{r_k}(\mathbb{E}[\text{TOPUP}_k](U_{-k})) = h_0^{1/2} \lambda_k^2$. By (3.8), the proof of (C.6) also implies

$$(C.8) \quad L^{(\text{ideal})} \leq C_1^{(\text{TOPUP})} R^{(\text{ideal})} \quad \text{with} \quad R^{(\text{ideal})} = \max_{1 \leq k \leq K} R_k^{(\text{ideal})}$$

in an event Ω_4 with $\bar{\mathbb{P}}(\Omega_4) \leq 5^{-1} \sum_{k=1}^K e^{-d_k}$, where $R_k^{(\text{ideal})}$ is as in (3.8).

It follows from (B.1), (C.4), (C.5) and (C.6) that in the event $\bigcap_{j=0}^4 \Omega_j$

$$(C.9) \quad L_k^{(m+1)} \leq L_k^{(\text{ideal})} + 24(K-1)L^{(m)} \max_{j,k,h} \|\Delta_{j,k,h}\|_{k,S,S} / \lambda_k^2 \leq L_k^{(\text{ideal})} + \rho L^{(m)},$$

which implies by induction

$$(C.10) \quad L^{(m+1)} \leq (1 + \dots + \rho^m) L^{(\text{ideal})} + \rho^{m+1} L^{(0)},$$

and then the conclusions follows from (C.1) and (C.8). Again, we divide the rest of the proof into 4 steps to prove (C.7) for $j = 1, 2, 3$ and (C.3).

Step 1. We prove (C.7) for the $\Delta_1(\tilde{U}_2, \tilde{U}_4)$. By Lemma G.1 (ii), there exist $U_2^{(\ell)}, U_2^{(\ell')} \in \mathbb{R}^{d_2 \times r_2}$, $1 \leq \ell, \ell' \leq N_{d_2 r_2, 1/8} := 17^{d_2 r_2}$, such that $\|U_2^{(\ell)}\|_S \leq 1$, $\|U_2^{(\ell')}\|_S \leq 1$ and

$$\|\Delta_1\|_{1,S,S} = \|\Delta_{1,1,h}\|_{1,S,S} \leq 2 \max_{\ell, \ell' \leq N_{d_2 r_2, 1/8}} \left\| \sum_{t=h+1}^T \frac{\text{mat}_1(A_1 F_{t-h} A_2^\top U_2^{(\ell)} \otimes E_t U_2^{(\ell')})}{T-h} \right\|_S.$$

We apply the Gaussian concentration inequality to the right-hand side above. Elementary calculation shows that

$$\begin{aligned} & \left\| \left\| \sum_{t=h+1}^T \text{mat}_1(A_1 F_{t-h} A_2^\top U_2^{(\ell)} \otimes E_t U_2^{(\ell')}) \right\|_S - \left\| \sum_{t=h+1}^T \text{mat}_1(A_1 F_{t-h} A_2^\top U_2^{(\ell)} \otimes E_t^* U_2^{(\ell')}) \right\|_S \right\| \\ & \leq \left\| \sum_{t=h+1}^T \text{mat}_1(A_1 F_{t-h} A_2^\top U_2^{(\ell)} \otimes (E_t - E_t^*) U_2^{(\ell')}) \right\|_S \\ & \leq \left\| (\text{mat}_1(A_1 F_1 A_2^\top \otimes I_{d_1}), \dots, \text{mat}_1(A_1 F_{T-h} A_2^\top \otimes I_{d_1})) \begin{pmatrix} U_2^{(\ell)} \odot I_{d_1} \odot (E_{h+1} - E_{h+1}^*) U_2^{(\ell')} \\ \vdots \\ U_2^{(\ell)} \odot I_{d_1} \odot (E_T - E_T^*) U_2^{(\ell')} \end{pmatrix} \right\|_S \\ & \leq \sqrt{T} \|\Theta_{1,0}^*\|_S^{1/2} \|U_2^{(\ell)}\|_S \|U_2^{(\ell')}\|_S \left\| \begin{pmatrix} E_{h+1} - E_{h+1}^* \\ \vdots \\ E_T - E_T^* \end{pmatrix} \right\|_F. \end{aligned}$$

That is, $\left\| \sum_{t=h+1}^T \text{mat}_1(A_1 F_{t-h} A_2^\top U_2^{(\ell)} \otimes E_t U_2^{(\ell')}) \right\|_S$ is a $\sqrt{T} \|\Theta_{1,0}^*\|_S^{1/2}$ Lipschitz function in (E_1, \dots, E_T) . Employing similar arguments in the proof of Theorem 1 in [Chen, Yang and Zhang \(2019\)](#), we have

$$\mathbb{E} \left\| \sum_{t=h+1}^T \frac{\text{mat}_1(A_1 F_{t-h} A_2^\top U_2^{(\ell)} \otimes E_t U_2^{(\ell')})}{T-h} \right\|_S \leq \frac{\sigma(2T)^{1/2}(\sqrt{d_1} + \sqrt{d_1 r_2^2})}{T-h} \|\Theta_{1,0}^*\|_S^{1/2}$$

Then, by Gaussian concentration inequalities for Lipschitz functions,

$$\begin{aligned} & \mathbb{P} \left(\left\| \sum_{t=h+1}^T \frac{\text{mat}_1(A_1 F_{t-h} A_2^\top U_2^{(\ell)} \otimes E_t U_2^{(\ell')})}{T-h} \right\|_S - \frac{\sigma(2T)^{1/2}(\sqrt{d_1} + \sqrt{d_1 r_2^2})}{T-h} \|\Theta_{1,0}^*\|_S^{1/2} \geq \frac{\sigma\sqrt{T}}{T-h} \|\Theta_{1,0}^*\|_S^{1/2} x \right) \\ & \leq 2e^{-\frac{x^2}{2}}. \end{aligned}$$

Hence,

$$\begin{aligned} & \mathbb{P} \left(\|\Delta_1\|_{1,S,S}/2 \geq \frac{\sigma(2T)^{1/2}(\sqrt{d_1} + \sqrt{d_1 r_2^2})}{T-h} \|\Theta_{1,0}^*\|_S^{1/2} + \frac{\sigma\sqrt{T}}{T-h} \|\Theta_{1,0}^*\|_S^{1/2} x \right) \\ & \leq 2N_{d_2 r_2, 1/8}^2 e^{-\frac{x^2}{2}}. \end{aligned}$$

As $T \geq 4h_0$, this implies with $x = \sqrt{d_2 r_2}$ that in an event with at least probability $1 - e^{-d_2/5}$,

$$\begin{aligned} \|\Delta_1\|_{1,S,S} & \leq \frac{C_{1,K}^{(\text{iter})} \sigma T^{-1/2} \|\Theta_{1,0}^*\|_S^{1/2} (\sqrt{d_1} r_2 + \sqrt{d_2 r_2})}{24(K-1)} \\ & \leq \frac{C_{1,K}^{(\text{iter})} (\lambda_1^2 R_1^{(\text{ideal})}) + \sigma T^{-1/2} \|\Theta_{1,0}^*\|_S^{1/2} \sqrt{d_2 r_2}}{24} \end{aligned}$$

with the $R_k^{(\text{ideal})}$ in (3.8) and a constant $C_{1,K}^{(\text{iter})}$ depending on K only. In this event, (3.10) gives $36(\lambda_k)^{-2} \|\Delta_{1,k,h}\|_{k,S,S} \leq C_{1,K}^{(\text{iter})} (R_k^{(\text{ideal})} + R_k^{(\text{add})}) \leq \rho$. Thus, (C.7) holds for $\Delta_1^*(\tilde{M}_2)$.

Step 2. Note that

$$\|\Delta_2\|_{1,S,S} = \max_{\substack{\tilde{U}_2 \in \mathbb{R}^{d_2 \times r_2}, \tilde{U}_4 \in \mathbb{R}^{d_2 \times r_2}, \\ \|\tilde{U}_2\|_S \leq 1, \|\tilde{U}_4\|_S \leq 1}} \left\| \sum_{t=h+1}^T \frac{\text{mat}_1(E_{t-h} \tilde{U}_2 \otimes U_1^\top A_1 F_t A_2^\top \tilde{U}_2)}{T-h} \right\|_S.$$

Then, inequality (C.7) for $\Delta_2(\tilde{U}_2, \tilde{U}_4)$ follow from the same argument as the above step.

Step 3. Now we prove (C.7) for the $\Delta_3(\tilde{U}_2, \tilde{U}_4)$. We split the sum into two terms over the index sets, $S_1 = \{(h, 2h] \cup (3h, 4h] \cup \dots\} \cap (h, T]$ and its complement S_2 in $(h, T]$, so that $\{E_{t-h}, t \in S_a\}$ is independent of $\{E_t, t \in S_a\}$ for each $a = 1, 2$. Let $n_a = |S_a|$.

By Lemma G.1 (ii), we can find $U_2^{(\ell)}, U_2^{(\ell')} \in \mathbb{R}^{d_2 \times r_2}$, $1 \leq \ell, \ell' \leq N_{d_2 r_2, 1/8}$ such that $\|U_2^{(\ell)}\|_S \leq 1$, $\|U_2^{(\ell')}\|_S \leq 1$. In this case,

$$(C.11) \quad \|\Delta_3\|_{1,S,S} = \|\Delta_{3,1,h}\|_{1,S,S} \leq 2 \max_{1 \leq \ell, \ell' \leq N_{d_2 r_2, 1/8}} \left\| \sum_{t=h+1}^T \frac{\text{mat}_1(E_{t-h} U_2^{(\ell)} \otimes E_t U_2^{(\ell')})}{T-h} \right\|_S.$$

Define $G_a = (E_{t-h} U_2^{(\ell)}, t \in S_a)$ and $H_a = (E_t U_2^{(\ell')}, t \in S_a)$. Then, G_a, H_a are two independent Gaussian matrices. By Lemma G.2(ii), for any $x > 0$,

$$\begin{aligned} \mathbb{P} \left(\left\| \sum_{t \in S_a} \text{mat}_1(E_{t-h} U_2^{(\ell)} \otimes E_t U_2^{(\ell')}) \right\|_S \geq d_1 \sqrt{r_2} + 2r_2 \sqrt{d_1 n_a} + x^2 + \sqrt{n_a} x + 3\sqrt{d_1 r_2 x} \right) \\ \leq 2e^{-x^2/2}. \end{aligned}$$

As in the derivation of $\|\Delta_3^*\|_{1,S,S}$ in the proof of Theorem 3.2, we have, with $x = \sqrt{d_2 r_2}$ and some constant $C_{1,K}^{(\text{iter})}$ depending on K only,

$$\mathbb{P} \left(\|\Delta_{3,1,h}\|_{1,S,S} \geq \frac{C_{1,K}^{(\text{iter})} \sigma^2}{24} \left(\frac{r_2 \sqrt{d_1} + \sqrt{d_2 r_2}}{T^{1/2}} + \frac{d_1 \sqrt{r_2} + d_2 r_2 + r_2 \sqrt{d_1 d_2}}{T} \right) \right) \leq e^{-d_2/5}.$$

This yields (C.7) for $\Delta_3(\tilde{U}_2, \tilde{U}_4)$ as in the end of Step 1 for $\Delta_1(\tilde{U}_2, \tilde{U}_4)$.

Step 4. Next, we consider the r_1 -th largest singular value of $\sigma_{r_1}(\overline{\mathbb{E}}[\text{TOPUP}_1](\hat{U}_2^{(m)}))$ in the event $\cap_{j=0}^4 \Omega_j$. By definition, the left singular subspace of $\overline{\mathbb{E}}[\text{TOPUP}_1](\hat{U}_2^{(m)})$ is U_1 . Then,

$$\begin{aligned} & \sigma_{r_1}(\overline{\mathbb{E}}[\text{TOPUP}_1](\hat{U}_2^{(m)})) \\ &= \sigma_{r_1} \left(\text{mat}_1 \left(\sum_{t=h+1}^T \frac{A_1 F_{t-h} A_2^\top \hat{U}_2^{(m)} \otimes A_1 F_t A_2^\top \hat{U}_2^{(m)}}{T-h}, h = 1, \dots, h_0 \right) \right) \\ &= \sigma_{r_1} \left(\text{mat}_1 \left(\sum_{t=h+1}^T \frac{A_1 F_{t-h} A_2^\top U_2 U_2^\top \hat{U}_2^{(m)} \otimes A_1 F_t A_2^\top U_2 U_2^\top \hat{U}_2^{(m)}}{T-h}, h = 1, \dots, h_0 \right) \right) \\ &= \sigma_{r_1} \left(\text{mat}_1 \left(\sum_{t=h+1}^T \frac{A_1 F_{t-h} A_2^\top U_2 \otimes A_1 F_t A_2^\top U_2}{T-h}, h = 1, \dots, h_0 \right) \cdot (U_2^\top \hat{U}_2^{(m)} \odot I_{d_1} \odot U_2^\top \hat{U}_2^{(m)} \odot I_{h_0}) \right) \end{aligned}$$

$$\begin{aligned}
&\geq \sigma_{r_1} \left(\text{mat}_1 \left(\sum_{t=h+1}^T \frac{A_1 F_{t-h} A_2^\top U_2 \otimes A_1 F_t A_2^\top U_2}{T-h}, h=1, \dots, h_0 \right) \right) \cdot \sigma_{\min} \left(U_2^\top \hat{U}_2^{(m)} \odot I_{d_1} \odot U_2^\top \hat{U}_2^{(m)} \odot I_{h_0} \right) \\
&\geq \sigma_{r_1} \left(\text{mat}_1 \left(\sum_{t=h+1}^T \frac{A_1 F_{t-h} A_2^\top \otimes A_1 F_t A_2^\top}{T-h}, h=1, \dots, h_0 \right) \right) \cdot \sigma_{\min} \left(U_2^\top \hat{U}_2^{(m)} \right) \cdot \sigma_{\min} \left(U_2^\top \hat{U}_2^{(m)} \right) \\
&\geq \sqrt{h_0} \lambda_1^2 (1 - L^{(m)2}).
\end{aligned}$$

The last step follows from the definitions in (2.1) and (B.1). If $L^{(m)} \leq 1/2$, then

$$\sigma_{r_1}(\mathbb{E}[\text{TOPUP}_1](\hat{U}_2^{(m)})) \geq \sqrt{h_0} \lambda_1^2 / 2.$$

By (3.5), $\lambda_1^2 h_0^{1/2} = \sigma_{r_1}(\text{mat}_1(\Theta_{1,1:h_0})) = \sigma_{r_1}(\text{mat}_1(\mathbb{E}[\text{TOPUP}_1](U_2)))$. We prove this condition in the event $\cap_{j=0}^4 \Omega_j$. By (3.10) and (C.1), $L^{(0)} \leq C_1^{(\text{TOPUP})} R^{(0)} \leq 1/2$. By induction, given $L^{(m)} \leq 1/2$, (C.3) holds for the same m . Applying (C.1), (C.8) and (C.10),

$$\begin{aligned}
L^{(m+1)} &\leq C_1^{(\text{TOPUP})} \{2(1 + \dots + \rho^m) R^{(\text{ideal})} + \rho^{m+1} R^{(0)}\} \leq C_1^{(\text{TOPUP})} 2(1 - \rho)^{-1} R^{(0)} \\
&\leq 1/2.
\end{aligned}$$

This completes the induction and the proof of the entire theorem.

APPENDIX D: PROOF OF THEOREM 3.4

Without loss of generality, we can assume $\sigma = 1$. Within the probability space (3.33) $\mathcal{P}(T, d_1, \dots, d_K, \lambda)$, we study a specific model with $f_{2t} = f_{2t-1}$ for all $1 \leq t \leq \lfloor T/2 \rfloor$. Taking an average for each \mathcal{X}_{2t} and \mathcal{X}_{2t-1} , $1 \leq t \leq \lfloor T/2 \rfloor$, we reduce by sufficiency the model to

$$\begin{aligned}
\tilde{\mathcal{P}}(\lfloor T/2 \rfloor, d_1, \dots, d_K, \lambda) &= \left\{ \mathcal{X}_1, \dots, \mathcal{X}_{\lfloor T/2 \rfloor} : \mathcal{X}_t = \lambda \tilde{f}_t \times_1 a_1 \times_2 \dots \times_K a_K + \tilde{\mathcal{E}}_t, \text{ with } a_k \in \mathbb{R}^{d_k}, \right. \\
&\quad \left\| a_k \right\|_2 = 1, 1 \leq k \leq K, \tilde{f}_t \stackrel{i.i.d.}{\sim} N(0, 1/2), \{\tilde{f}_t\}_{t=1}^{\lfloor T/2 \rfloor} \text{ independent of } \{\tilde{\mathcal{E}}_t\}_{t=1}^{\lfloor T/2 \rfloor}, \\
&\quad \left. \tilde{\mathcal{E}}_{t, j_1, \dots, j_K} \stackrel{i.i.d.}{\sim} N(0, 1/2) \text{ for all } 1 \leq t \leq \lfloor T/2 \rfloor, 1 \leq j_k \leq d_k, 1 \leq k \leq K \right\}.
\end{aligned}$$

For notation convenience, in the following of this section, we study the probability space

(D.1)

$$\begin{aligned}
\tilde{\mathcal{P}}(T, d_1, \dots, d_K, \lambda) &= \left\{ \mathcal{X}_1, \dots, \mathcal{X}_T : \mathcal{X}_t = \lambda f_t \times_1 a_1 \times_2 \dots \times_K a_K + \mathcal{E}_t, \text{ with } a_k \in \mathbb{R}^{d_k}, \right. \\
&\quad \left\| a_k \right\|_2 = 1, 1 \leq k \leq K, f_t \stackrel{i.i.d.}{\sim} N(0, 1), \{f_t\}_{t=1}^T \text{ independent of } \{\mathcal{E}_t\}_{t=1}^T, \\
&\quad \left. \mathcal{E}_{t, j_1, \dots, j_K} \stackrel{i.i.d.}{\sim} N(0, 1) \text{ for all } 1 \leq t \leq T, 1 \leq j_k \leq d_k, 1 \leq k \leq K \right\}.
\end{aligned}$$

We first introduce some additional notation. For any probability distributions \mathbb{P} and \mathbb{Q} , define total variation distance as $\text{TV}(\mathbb{P}, \mathbb{Q}) = \sup_B |\mathbb{P}(B) - \mathbb{Q}(B)|$. We also write $\text{TV}(p, q)$ if p, q are the densities of \mathbb{P} and \mathbb{Q} , respectively. Define

$$\tau_N = \frac{\kappa}{N}, \quad \eta_N = \frac{\kappa}{(K+2)^4 N (\log N)^2},$$

where κ is the size of the clique. For any $\mu \in \mathbb{R}$, denote ϕ_μ as the density function of the $N(\mu, 1)$ distribution, and let

$$\bar{\phi}_\mu = \frac{1}{2}(\phi_\mu + \phi_{-\mu})$$

be the density function of the Gaussian mixture $\frac{1}{2}N(\mu, 1) + \frac{1}{2}N(-\mu, 1)$. We also define $\tilde{\Xi}_0$ as a truncated normal distribution by the $N(0, 1)$ distribution restricted on the interval $[-(K+2)\sqrt{\log N}, (K+2)\sqrt{\log N}]$. For any $|\mu| \leq (K+2)\sqrt{\eta_N \log N}$, define two distributions $\mathcal{F}_{\mu,0}$ and $\mathcal{F}_{\mu,1}$ with density functions

$$\begin{aligned} h_{\mu,0}(x) &= J_0(\phi_0(x) - \tau_N^{-1}[\bar{\phi}_\mu(x) - \phi_0(x)])\mathbf{1}_{\{|x| \leq (K+2)\sqrt{\log N}\}}, \\ h_{\mu,1}(x) &= J_1(\phi_0(x) + \tau_N^{-1}[\bar{\phi}_\mu(x) - \phi_0(x)])\mathbf{1}_{\{|x| \leq (K+2)\sqrt{\log N}\}}, \end{aligned}$$

where J_0, J_1 are normalizing constants.

Suppose we have a collection of estimators $\hat{\mathbf{a}} = (\hat{a}_1, \dots, \hat{a}_K)$ with $\hat{a}_j = \hat{a}_j(\mathcal{X}_1, \dots, \mathcal{X}_T)$ being the estimated factor loading a_j ($1 \leq j \leq K$). Our main technique is based on a reduction scheme which maps any order $K+1$ graph adjacent tensor $\mathcal{A} \in \{0, 1\}^{N^{\otimes(K+1)}}$ with dimension $N^{\otimes(K+1)} = N \times N \times \dots \times N$ (multiply N by $K+1$ times), $N \geq 2T$, and $\hat{\mathbf{a}}$, to a test for the Hypergraphic Planted Clique detection problem (3.30). The specific technique was developed by [Ma and Wu \(2015\)](#), [Gao, Ma and Zhou \(2017\)](#). We refer the readers to [Wang, Berthet and Samworth \(2016\)](#), [Cai, Liang and Rakhlin \(2017\)](#) for other related methods. We provide a detailed description of the mapping as follows.

(1) (Initialization). Generate i.i.d. random variable $\xi_1, \dots, \xi_{2T} \sim \tilde{\Xi}_0$. Set

$$(D.2) \quad \mu_t = \eta_N^{1/2} \xi_t, \quad t = 1, \dots, 2T.$$

(2) (Gaussianization). Generate two order $K+1$ tensors $\mathcal{B}_0, \mathcal{B}_1 \in \mathbb{R}^{2T^{\otimes(K+1)}}$, where conditioning on the μ_t 's, all the entries are mutually independent satisfying

$$(D.3) \quad \mathcal{L}((\mathcal{B}_0)_{t,j_1,\dots,j_K} | \mu_t) = \mathcal{F}_{\mu_t,0} \quad \text{and} \quad \mathcal{L}((\mathcal{B}_1)_{t,j_1,\dots,j_K} | \mu_t) = \mathcal{F}_{\mu_t,1}.$$

Let $\mathcal{A}_0 \in \{0, 1\}^{2T^{\otimes(K+1)}}$ be the lower-left corner block of the order $K+1$ tensor \mathcal{A} . Generate an order $K+1$ tensor $\mathcal{X} = (\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_{2T}) \in \mathbb{R}^{2T \times d_1 \times d_2 \times \dots \times d_K}$, where for each $t \leq 2T$, if $1 \leq j_1, \dots, j_K \leq 2T$, set

$$(D.4) \quad \mathcal{X}_{t,j_1,\dots,j_K} = (1 - (\mathcal{A}_0)_{t,j_1,\dots,j_K})(\mathcal{B}_0)_{t,j_1,\dots,j_K} + (\mathcal{A}_0)_{t,j_1,\dots,j_K}(\mathcal{B}_1)_{t,j_1,\dots,j_K};$$

otherwise, let $\mathcal{X}_{t,j_1,\dots,j_K}$ be an independent copy from $N(0, 1)$.

(3) (Test Construction). Let $\hat{a}_k = \hat{a}_k(\mathcal{X}_1, \dots, \mathcal{X}_T)$ be the estimator of the factor loading a_k , for all $1 \leq k \leq K$, by treating $\mathcal{X}_1, \dots, \mathcal{X}_T$ as data. All \hat{a}_k are normalized to be a unit vector. We reject H_0^G if

$$(D.5) \quad \left(\frac{1}{T} \sum_{t=T+1}^{2T} \mathcal{X}_t \otimes \mathcal{X}_t \right) \times_{k=1}^{2K} \hat{a}_k^\top \geq 1 + \frac{1}{2} \left(\frac{\kappa}{2} \right)^K \eta_N,$$

with $\hat{a}_{k+K} := \hat{a}_k$ for all $1 \leq k \leq K$.

D.1. Lemmas. To prove Theorem 3.4, we need to state lemmas which characterize the distribution of \mathcal{X}_t under H_0^G and H_1^G . Let $\mathcal{L}(\{\mathcal{X}_t\}_{t=1}^{2T})$ be the joint distribution of $\mathcal{X}_1, \dots, \mathcal{X}_{2T}$. Under probability space (D.1), we denote the distribution of \mathcal{X}_t as $\mathbb{P}_{\lambda, \mathbf{a}}$ with $\mathbf{a} = (a_1, \dots, a_K)$, and the joint distribution of $(\mathcal{X}_1, \dots, \mathcal{X}_{2T})$ as $\mathbb{P}_{\lambda, \mathbf{a}}^{2T}$. Note that $\mathbb{P}_{\lambda, \mathbf{a}}$ is a tensor (array) normal distribution with covariance tensor $\lambda \times_{k=1}^{2K} a_k + \mathcal{I}_{d_1 \times \dots \times d_K \times d_1 \times \dots \times d_K}$, where $a_{K+k} = a_k$ for $1 \leq k \leq K$, $\mathcal{I}_{j_1, \dots, j_K, j_1, \dots, j_K} = 1$ and 0 elsewhere. If $\lambda = 0$, all entries of \mathcal{X}_t are i.i.d. $N(0, 1)$, thus we denote $\mathcal{L}(\mathcal{X}_t)$ as \mathbb{P}_0 . We also write $\mathcal{L}(\mathcal{X}|\beta)$ as the conditional distribution of $\mathcal{X}|\beta$, and $\int \mathcal{L}(\mathcal{X}|\beta) d\xi(\beta)$ as the marginal distribution of \mathcal{X} after integrating β out.

LEMMA D.1. *There exists some absolute constant $C > 0$, such that for any integers $K \geq 1$, $\kappa < N$, $N \geq 3$, and for all $|\mu| \leq (K + 2)\sqrt{\eta_N \log N}$,*

$$\text{TV}(g_{\mu,0}, \phi_0) \leq CN^{-K-2} \quad \text{and} \quad \text{TV}(g_{\mu,1}, \bar{\phi}_\mu) \leq CN^{-K-2},$$

where $g_{\mu,0} = \frac{1}{2}(h_{\mu,0} + h_{\mu,1})$ and $g_{\mu,1} = \tau_N h_{\mu,1} + (1 - \tau_N)\frac{1}{2}(h_{\mu,0} + h_{\mu,1})$.

LEMMA D.2. *Suppose $\mathcal{A} \sim \mathcal{G}_{K+1}(N, 1/2)$. There exists some constant $C_K > 0$ depending on K only, such that*

$$\text{TV}(\mathcal{L}(\{\mathcal{X}_t\}_{t=1}^{2T}, \mathbb{P}_0^{2T}) \leq C_K N^{-1}.$$

The proofs of Lemma D.1 and D.2 are analogous to Lemma 7.1 and 7.2 in [Gao, Ma and Zhou \(2017\)](#), thus are skipped here.

LEMMA D.3. *Suppose $\mathcal{A} \sim \mathcal{G}_{K+1}(N, 1/2, \kappa)$. There exists a distribution π supported on the set*

$$\{(\lambda, \mathbf{a}) : \|\mathbf{a}_k\|_2 = 1, |\text{supp}(\mathbf{a}_k)| \leq 3\kappa/2, 1 \leq k \leq K, \eta_N^{1/2}(\kappa/2)^{K/2} \leq \lambda \leq \eta_N^{1/2}(3\kappa/2)^{K/2}\},$$

such that for some positive constants C_{1K}, C_{2K} depending on K only,

$$\text{TV}(\mathcal{L}(\{\mathcal{X}_t\}_{t=1}^{2T}, \mathbb{P}_\pi) \leq C_{1K} \cdot \kappa \left(\frac{2T}{N}\right)^\kappa + \frac{C_{2K}}{N} + \frac{4(K+1)T}{N},$$

where $\mathbb{P}_\pi = \int \mathbb{P}_{\lambda, \mathbf{a}}^{2T} d\pi(\lambda, \mathbf{a})$.

PROOF. Let ξ be $N(0, \eta_N)$, and $\bar{\xi}$ be a truncated normal distribution obtained by restricting ξ on the set $[-(K+2)\sqrt{\eta_N \log N}, (K+2)\sqrt{\eta_N \log N}]$. Then the μ_i 's in (D.2) are i.i.d. following $\bar{\xi}$. Elementary calculation shows that $\int \phi_0(x) d\xi(\mu) = \phi_0(x)$ gives the density function of $N(0, 1)$, and $\int \bar{\phi}_\mu(x) d\xi(\mu)$ is the density function of $N(0, 1 + \eta_N)$.

We first consider the case $d_1 = d_2 = \dots = d_K = 2T$. Let $(\alpha_1, \dots, \alpha_{2T})$ be the 0-1 indicators of the first tensor mode of \mathcal{A}_0 whether the corresponding vertices belong to the planted clique or not. Similarly, define $(\beta_{k1}, \dots, \beta_{k,2T})$ as the corresponding indicators of the $(k+1)$ -th tensor mode of \mathcal{A}_0 , for all $1 \leq k \leq K$. Let $(\tilde{\alpha}_1, \dots, \tilde{\alpha}_{2T})$ and $(\tilde{\beta}_{k1}, \dots, \tilde{\beta}_{k,2T})$, $1 \leq k \leq K$, as i.i.d. Bernoulli random variables with mean $\tau_N = \kappa/N$. Define a new tensor $\tilde{\mathcal{A}}_0$ with $\tilde{\mathcal{A}}_{t,j_1, \dots, j_K} = 1$ if $\tilde{\alpha}_t = \tilde{\beta}_{1,j_1} = \dots = \tilde{\beta}_{K,j_K} = 1$ and is an independent instantiation of the Bernoulli(1/2) distribution otherwise. Define $\tilde{\mathcal{X}}$ as

$$\tilde{\mathcal{X}}_{t,j_1, \dots, j_K} = (1 - (\tilde{\mathcal{A}}_0)_{t,j_1, \dots, j_K})(\mathcal{B}_0)_{t,j_1, \dots, j_K} + (\tilde{\mathcal{A}}_0)_{t,j_1, \dots, j_K}(\mathcal{B}_1)_{t,j_1, \dots, j_K}.$$

By Theorem 4 in [Diaconis and Freedman \(1980\)](#) and the data-processing inequality, we have

$$\text{TV}(\mathcal{L}(\mathcal{X}), \mathcal{L}(\tilde{\mathcal{X}})) \leq \text{TV}(\mathcal{L}(\alpha, \beta_1, \dots, \beta_K), \mathcal{L}(\tilde{\alpha}, \tilde{\beta}_1, \dots, \tilde{\beta}_K)) \leq \frac{4(K+1)T}{N},$$

where $\alpha = (\alpha_1, \dots, \alpha_{2T})$, $\beta_k = (\beta_{k1}, \dots, \beta_{k,2T})$, $1 \leq k \leq K$, and $\tilde{\alpha}, \tilde{\beta}_k$ are similarly defined. Note that, conditioning on μ_t and $\tilde{\beta}_{1,j_k} = 0$ for some $1 \leq k \leq K$, $\tilde{\mathcal{X}}_{t,j_1, \dots, j_K} \sim g_{\mu_t, 0}$. And conditioning on μ_t and $\tilde{\beta}_{1,j_1} = \dots = \tilde{\beta}_{K,j_K} = 1$, $\tilde{\mathcal{X}}_{t,j_1, \dots, j_K} \sim g_{\mu_t, 1}$.

Next, define $\bar{\mathcal{X}}$ with entries

$$\bar{\mathcal{X}}_{t,j_1, \dots, j_K} | (\tilde{\beta}_{1,j_k} = 0, \text{ for some } 1 \leq k \leq K, \mu_t) \sim \phi_0,$$

$$\bar{\mathcal{X}}_{t,j_1, \dots, j_K} | (\tilde{\beta}_{1,j_1} = \dots = \tilde{\beta}_{K,j_K} = 1, \mu_t) \sim \bar{\phi}_{\mu_t}.$$

By Lemma D.1 and Lemma 7 in Ma and Wu (2015), we have, uniformly over $\max_t |\mu_t| \leq (K+2)\sqrt{\eta_N \log N}$,

$$\begin{aligned} & \text{TV}(\mathcal{L}(\tilde{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K, \mu_t), \mathcal{L}(\bar{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K, \mu_t)) \\ & \leq \sum_{t=1}^{2T} \sum_{j_1=1}^{d_1} \cdots \sum_{j_K=1}^{d_K} \text{TV}(\mathcal{L}(\tilde{\mathcal{X}}_{t,j_1, \dots, j_K}|\tilde{\beta}_{1,j_1}, \dots, \tilde{\beta}_{K,j_K}, \mu_t), \mathcal{L}(\bar{\mathcal{X}}_{t,j_1, \dots, j_K}|\tilde{\beta}_{1,j_1}, \dots, \tilde{\beta}_{K,j_K}, \mu_t)) \\ & \leq \frac{C_K}{N}, \end{aligned}$$

for some constant $C_K > 0$. Let $\int \mathcal{L}(\tilde{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K, \mu) d\bar{\xi}(\mu)$ (resp. $\int \mathcal{L}(\tilde{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K, \mu) d\xi(\mu)$) be the conditional distribution of $\tilde{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K$ if the elements of $\mu = (\mu_1, \dots, \mu_{2T})$ are i.i.d. following $\bar{\xi}$ (resp. ξ). And $\int \mathcal{L}(\bar{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K, \mu) d\bar{\xi}(\mu)$, $\int \mathcal{L}(\bar{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K, \mu) d\xi(\mu)$ are similarly defined. Note that

$$\text{TV}(\xi, \bar{\xi}) = \int_{|\mu| > (K+2)\sqrt{\eta_N \log N}} d\xi(\mu) = \int_{|x| > (K+2)\sqrt{\log N}} \phi_0(x) dx \leq CN^{-K-3}.$$

Then, we can obtain

$$\begin{aligned} & \text{TV}\left(\int \mathcal{L}(\tilde{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K, \mu) d\bar{\xi}(\mu), \int \mathcal{L}(\bar{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K, \mu) d\xi(\mu)\right) \\ & \leq \text{TV}\left(\int \mathcal{L}(\tilde{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K, \mu) d\bar{\xi}(\mu), \int \mathcal{L}(\bar{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K, \mu) d\bar{\xi}(\mu)\right) \\ & \quad + \text{TV}\left(\int \mathcal{L}(\bar{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K, \mu) d\bar{\xi}(\mu), \int \mathcal{L}(\bar{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K, \mu) d\xi(\mu)\right) \\ & \leq \sup_{\max_t |\mu_t| \leq (K+2)\sqrt{\eta_N \log N}} \text{TV}(\mathcal{L}(\tilde{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K, \mu_t), \mathcal{L}(\bar{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K, \mu_t)) + C_0(2T)^K \text{TV}(\bar{\xi}, \xi) \\ & \leq C_K N^{-1}. \end{aligned}$$

Define a set

$$S_T := \{(j_1, j_2, \dots, j_K) : \tilde{\beta}_{1,j_1} = \dots = \tilde{\beta}_{K,j_K} = 1, 1 \leq j_k \leq d_k, 1 \leq k \leq K\}.$$

Then, for each given $(\tilde{\beta}_1, \dots, \tilde{\beta}_K)$, we can define $s_k = \sum_{j_k \in S_T} \tilde{\beta}_{k,j_k} = \sum_{j_k \in S_T} \tilde{\beta}_{k,j_k}^2$, $a_k = s_k^{-1/2}(\tilde{\beta}_{k,j_k} \mathbf{1}_{\{j_k \in S_T\}})$, for all $1 \leq k \leq K$, and $\lambda = \eta_N^{1/2} \prod_{k=1}^K s_k^{1/2}$. Obviously, there exists one-to-one identification between $(a_1, \dots, a_K, \lambda)$ and $(\tilde{\beta}_1, \dots, \tilde{\beta}_K)$. Note that $\int \mathcal{L}(\bar{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K, \mu) d\xi(\mu) = \mathbb{P}_{\lambda, \mathbf{a}}^{2T}$. As $\mathcal{L}(\tilde{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K) = \int \mathcal{L}(\tilde{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K, \mu) d\bar{\xi}(\mu)$, we have

$$\text{TV}(\mathcal{L}(\tilde{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K), \mathbb{P}_{\lambda, \mathbf{a}}^{2T}) \leq C_K N^{-1}.$$

Define an event $Q = \{\tilde{\beta}_1, \dots, \tilde{\beta}_K : |s_k - \kappa| \leq \kappa/2, 1 \leq k \leq K\}$. By Lemma D.4, $\mathbb{P}(Q^c) \leq K\kappa(2T/N)^\kappa$. Let $\tilde{\pi}$ be the joint distribution of (\mathbf{a}, λ) , and π be the distribution by restricting $\tilde{\pi}$ on $\{a_1(\tilde{\beta}_1), \dots, a_K(\tilde{\beta}_K), \lambda(\tilde{\beta}_1, \dots, \tilde{\beta}_K) : \tilde{\beta}_1, \dots, \tilde{\beta}_K \in Q\}$. It follows that $\text{TV}(\tilde{\pi}, \pi) \leq C\mathbb{P}(Q^c) \leq CK\kappa(2T/N)^\kappa$. As $\mathcal{L}(\bar{\mathcal{X}}|\tilde{\beta}_1, \dots, \tilde{\beta}_K) = \mathcal{L}(\bar{\mathcal{X}}|a_1, \dots, a_K, \lambda) = \mathcal{L}(\bar{\mathcal{X}}|\mathbf{a}, \lambda)$,

$$\begin{aligned} \text{TV}(\mathcal{L}(\bar{\mathcal{X}}), \mathbb{P}_\pi) & \leq \text{TV}(\mathcal{L}(\bar{\mathcal{X}}), \int \mathcal{L}(\bar{\mathcal{X}}|\mathbf{a}, \lambda) d\pi(\mathbf{a}, \lambda)) \\ & \quad + \text{TV}\left(\int \mathcal{L}(\bar{\mathcal{X}}|\mathbf{a}, \lambda) d\pi(\mathbf{a}, \lambda), \int \mathbb{P}_{\lambda, \mathbf{a}}^{2T} d\pi(\lambda, \mathbf{a})\right) \end{aligned}$$

$$\begin{aligned} &\leq \text{TV}(\tilde{\pi}, \pi) + \sup_{\mathbf{a}, \lambda} \text{TV}(\mathcal{L}(\tilde{\mathcal{X}}|\mathbf{a}, \lambda), \mathbb{P}_{\lambda, \mathbf{a}}^{2T}) \\ &\leq CK\kappa \left(\frac{2T}{N}\right)^\kappa + \frac{C_K}{N}. \end{aligned}$$

Hence, we have

$$\text{TV}(\mathcal{L}(\{\mathcal{X}_t\}_{t=1}^{2T}), \int \mathbb{P}_{\lambda, \mathbf{a}}^{2T} d\pi(\lambda, \mathbf{a})) \leq C_{1K} \cdot \kappa \left(\frac{2T}{N}\right)^\kappa + \frac{C_{2K}}{N} + \frac{4(K+1)T}{N}.$$

When $d_k \geq 2T$, $1 \leq k \leq K$, we first use the above arguments to analyze the distribution of the first $2T$ coordinates. Then, as the remaining $2T$ coordinates are exact, the total variation bound is zero. \square

LEMMA D.4. *Let $s_k = \sum_{j_k \in S_T} \tilde{\beta}_{k, j_k} = \sum_{j_k \in S_T} \tilde{\beta}_{k, j_k}^2$, $d_k = 2T < N$ for all $1 \leq k \leq K$. Define an event $Q = \{\tilde{\beta}_1, \dots, \tilde{\beta}_K : |s_k - \kappa| \leq \kappa/2, 1 \leq k \leq K\}$, then*

$$\mathbb{P}(Q) \geq 1 - \frac{K(\kappa+1)}{2} \left(\frac{2T}{N}\right)^\kappa.$$

PROOF. Recall that $\kappa < \sqrt{N}$. Let $\mathcal{C}_k = \{\tilde{\beta}_{k, j_k} : j_k \in S_T\}$, $1 \leq k \leq K$. Then

$$\begin{aligned} \mathbb{P}(|\mathcal{C}_k| \leq \kappa/2) &\leq \frac{\sum_{i=0}^{\kappa/2} \binom{d_k}{i}}{\binom{N}{\kappa}} \leq \frac{\kappa+1}{2} \cdot \frac{\binom{d_k}{\kappa/2}}{\binom{N}{\kappa}} = \frac{\kappa+1}{2} \cdot \frac{d_k(d_k-1)\cdots(d_k-\kappa/2+1) \cdot \kappa!}{(\kappa/2)! \cdot N(N-1)\cdots(N-\kappa+1)} \\ &\leq \frac{\kappa+1}{2} \cdot \left(\frac{d_k}{N}\right)^\kappa. \end{aligned}$$

Therefore, by Bonferroni inequality, we have the desired result. \square

D.2. Proof of Theorem 3.4. Write $\hat{\Sigma} = \frac{1}{T} \sum_{t=T+1}^{2T} \mathcal{X}_t \otimes \mathcal{X}_t$. Then the test (D.5) can be rewritten as

$$\psi = \psi(\mathcal{X}_1, \dots, \mathcal{X}_{2T}) = \psi(\mathcal{A}, \mu, \mathcal{B}_0, \mathcal{B}_1) := \mathbf{1} \left\{ \hat{\Sigma} \times_{k=1}^{2K} \hat{a}_k^\top \geq 1 + \frac{1}{2} \left(\frac{\kappa}{2}\right)^K \eta_N \right\}.$$

Note that ψ is a test for the Hypergraphic Planted Clique detection problem (3.30). Recall the probability space (D.1). For any (λ, \mathbf{a}) in the support of π , we have

$$\mathbb{P}_{\lambda, \mathbf{a}}^T \in \tilde{\mathcal{P}}(T, d_1, \dots, d_K, \lambda)$$

with $\eta_N^{1/2} (\kappa/2)^{K/2} \leq \lambda \leq \eta_N^{1/2} (3\kappa/2)^{K/2}$.

We first bound the Type-I error of the test ψ . By Lemma D.2,

$$\mathbb{P}_{H_0^G}(\psi = 1) \leq \mathbb{P}_0^T(\psi = 1) + C_K N^{-1}.$$

Under \mathbb{P}_0^T , $\hat{\mathbf{a}} = (\hat{a}_1, \dots, \hat{a}_K)$ and $\hat{\Sigma}$ are independent. Conditioning on $\hat{\mathbf{a}}$, applying Bernstein's inequality, we have

$$\hat{\Sigma} \times_{k=1}^{2K} \hat{a}_k^\top = 1 + \frac{1}{T} \sum_{t=T+1}^{2T} \left(|\mathcal{X}_t \times_{k=1}^K \hat{a}_k^\top|^2 - |\text{vec}(\otimes_{k=1}^K \hat{a}_k)|^2 \right) > 1 + \frac{1}{2} \left(\frac{\kappa}{2}\right)^K \eta_N,$$

with probability at most $\exp\left(-\frac{C_K T \kappa^{K+2}}{N^2 (\log N)^4}\right)$. Integrating over $\hat{\mathbf{a}}$, we have

$$(D.6) \quad \mathbb{P}_{H_0^G}(\psi = 1) \leq \exp\left(-\frac{C_K T \kappa^{K+2}}{N^2 (\log N)^4}\right) + C_K N^{-1}.$$

Next, we bound the Type-II error. By Lemma D.3,

$$\mathbb{P}_{H_1^c}(\psi = 0) \leq \mathbb{P}_\pi(\psi = 0) + C_{1K} \cdot \kappa \left(\frac{2T}{N} \right)^\kappa + \frac{C_{2K}}{N} + \frac{4(K+1)T}{N}.$$

Recall that under the probability space (D.1),

$$\mathcal{X}_t = \lambda f_t \times_1 a_1 \times_2 \dots \times_K a_K + \mathcal{E}_t,$$

and $f_t \stackrel{i.i.d.}{\sim} N(0, 1)$ and the elements of \mathcal{E}_t follow that $\mathcal{E}_{t, j_1, \dots, j_K} \stackrel{i.i.d.}{\sim} N(0, 1)$ for all $1 \leq t \leq T, 1 \leq j_k \leq d_k, 1 \leq k \leq K$. Write $\mathcal{E}_t(\hat{\mathbf{a}}) = \mathcal{E}_t \times_{k=1}^K \hat{a}_k^\top$. Thus,

$$\begin{aligned} \hat{\Sigma} \times_{k=1}^{2K} \hat{a}_k^\top &= \lambda^2 \prod_{k=1}^K |\hat{a}_k^\top a_k|^2 \left(\frac{1}{T} \sum_{t=T+1}^{2T} f_t^2 \right) + \frac{2\lambda}{T} \sum_{t=T+1}^{2T} \prod_{k=1}^K (\hat{a}_k^\top a_k) \mathcal{E}_t(\hat{\mathbf{a}}) + \frac{1}{T} \sum_{t=T+1}^{2T} |\mathcal{E}_t(\hat{\mathbf{a}})|^2 \\ &= \lambda^2 + 1 + \lambda^2 \prod_{k=1}^K |\hat{a}_k^\top a_k|^2 \left(\frac{1}{T} \sum_{t=T+1}^{2T} f_t^2 - 1 \right) + \lambda^2 \left(\prod_{k=1}^K |\hat{a}_k^\top a_k|^2 - 1 \right) \\ &\quad + \frac{2\lambda}{T} \sum_{t=T+1}^{2T} \prod_{k=1}^K (\hat{a}_k^\top a_k) \mathcal{E}_t(\hat{\mathbf{a}}) + \left(\frac{1}{T} \sum_{t=T+1}^{2T} |\mathcal{E}_t(\hat{\mathbf{a}})|^2 - 1 \right). \end{aligned}$$

After rearrangement, we have

$$\begin{aligned} \left| \hat{\Sigma} \times_{k=1}^{2K} \hat{a}_k^\top - (\lambda^2 + 1) \right| &\leq \lambda^2 \left| \frac{1}{T} \sum_{t=T+1}^{2T} f_t^2 - 1 \right| + \lambda^2 \left(1 - \prod_{k=1}^K |\hat{a}_k^\top a_k|^2 \right) \\ &\quad + \left| \frac{2\lambda}{T} \sum_{t=T+1}^{2T} \mathcal{E}_t(\hat{\mathbf{a}}) \right| + \left| \frac{1}{T} \sum_{t=T+1}^{2T} |\mathcal{E}_t(\hat{\mathbf{a}})|^2 - 1 \right|. \end{aligned}$$

By Bernstein's inequality, we can obtain

$$\begin{aligned} \mathbb{P}_{\lambda, \mathbf{a}}^T \left(\left| \frac{1}{T} \sum_{t=T+1}^{2T} f_t^2 - 1 \right| + \left| \frac{1}{T} \sum_{t=T+1}^{2T} \mathcal{E}_t(\hat{\mathbf{a}}) \right| + \left| \frac{1}{T} \sum_{t=T+1}^{2T} |\mathcal{E}_t(\hat{\mathbf{a}})|^2 - 1 \right| \geq C \sqrt{\frac{\log T}{T}} \right) \\ \leq T^{-C'}. \end{aligned}$$

Note that

$$\lambda^2 \left(1 - \prod_{k=1}^K |\hat{a}_k^\top a_k|^2 \right) \geq \lambda^2 - \lambda^2 \max_k |\hat{a}_k^\top a_k|^2 = \lambda^2 \min_k \|P_{\hat{a}_k} - P_{a_k}\|_S^2.$$

Hence, as $\lambda^2 \geq \eta_N(\kappa/2)^K$, the Type-II error is upper bounded by

$$\begin{aligned} \text{(D.7)} \quad \mathbb{P}_{H_1^c}(\psi = 0) &\leq \mathbb{P}_{\lambda, \mathbf{a}}^T \left(\min_k \|P_{\hat{a}_k} - P_{a_k}\|_S^2 > \frac{1}{3} \right) + T^{-C'} + C_{1K} \cdot \kappa \left(\frac{2T}{N} \right)^\kappa \\ &\quad + \frac{C_{2K}}{N} + \frac{4(K+1)T}{N}. \end{aligned}$$

Combining (D.6) and (D.7), we have

$$\begin{aligned} \text{(D.8)} \quad \mathbb{P}_{H_0^c}(\psi = 1) + \mathbb{P}_{H_1^c}(\psi = 0) \\ \leq \mathbb{P}_{\lambda, \mathbf{a}}^T \left(\min_k \|P_{\hat{a}_k} - P_{a_k}\|_S^2 > \frac{1}{3} \right) + \exp \left(-\frac{C_K T \kappa^{K+2}}{N^2 (\log N)^4} \right) + C_K N^{-1} \\ + T^{-C'} + C_{1K} \cdot \kappa \left(\frac{2T}{N} \right)^\kappa + \frac{C_{2K}}{N} + \frac{4(K+1)T}{N}. \end{aligned}$$

Now, we are ready to prove Theorem 3.4. Consider the Hypergraphic Planted Clique detection problem (3.30) with $N = 20(K+1)T$, $\kappa = \lfloor N^{1/2-\delta} \rfloor$ and $\delta < 1/2 - 1/(K+2)$. Then we have

$$(D.9) \quad \frac{N(\log N)^5}{\kappa^{K+2}} \leq c_0,$$

for some sufficient small constant $c_0 > 0$. We can also obtain

$$\lambda^2 \leq \eta_N \left(\frac{3\kappa}{2} \right)^K \leq \frac{C_K d^{1/2-\delta(K+1)/K}}{T^{1/2}(\log T)^2}.$$

Obviously, (3.34) holds with $\vartheta = \delta(K+1)/K < (K+1)/(2K+4)$. On the contrary, suppose that the claim of Theorem 3.4 does not hold. It means that

$$(D.10) \quad \liminf_{T \rightarrow \infty} \sup_{\mathcal{X}_1, \dots, \mathcal{X}_T \in \mathcal{P}(T, d_1, \dots, d_K, \lambda)} \mathbb{P} \left(\min_{1 \leq k \leq K} \|P_{\hat{a}_k} - P_{a_k}\|_{\mathbb{S}}^2 > \frac{1}{3} \right) \leq \frac{1}{4}.$$

Substituting (D.9) and (D.10) into (D.8), we have

$$\begin{aligned} \mathbb{P}_{H_0^G}(\psi = 1) + \mathbb{P}_{H_1^G}(\psi = 0) &\leq \frac{1}{4}(1 + o(1)) + \frac{1}{5} + N^{-C_0\kappa} + C_K N^{-1} + T^{-C'} \\ &\quad + C_{1K} \cdot \kappa \left(\frac{1}{10K} \right)^\kappa + \frac{C_{2K}}{N}. \end{aligned}$$

It follows that

$$\limsup_{N \rightarrow \infty} (\mathbb{P}_{H_0^G}(\psi = 1) + \mathbb{P}_{H_1^G}(\psi = 0)) < \frac{1}{2},$$

which contradicts the Hypothesis I. We complete the proof.

APPENDIX E: PROOF OF THEOREM 3.5

Without loss of generality, we can assume $\sigma = 1$. Applying the same reduction in Appendix D, we map the probability space (3.33) to (D.1). Under probability space (D.1), we denote the distribution of \mathcal{X}_t as $\mathbb{P}_{\lambda, \mathbf{a}}$ with $\mathbf{a} = (a_1, \dots, a_K)$, and the joint distribution of $(\mathcal{X}_1, \dots, \mathcal{X}_T)$ as $\mathbb{P}_{\lambda, \mathbf{a}}^T$. We first present a lemma on the Kullback-Leibler divergence between data distributions generated by a special kind of tensor factor models.

LEMMA E.1. *For $i = 1, 2$, let $\mathbf{a}^{(i)} = (a_1^{(i)}, \dots, a_K^{(i)})$. Then the Kullback-Leibler divergence of $\mathbb{P}_{\lambda, \mathbf{a}^{(2)}}^T$ with respect to $\mathbb{P}_{\lambda, \mathbf{a}^{(1)}}^T$ is given by*

$$(E.1) \quad D(\mathbb{P}_{\lambda, \mathbf{a}^{(1)}}^T \| \mathbb{P}_{\lambda, \mathbf{a}^{(2)}}^T) = \frac{T\lambda^4}{2(1+\lambda^2)} \left(1 - \prod_{k=1}^K |a_k^{(1)\top} a_k^{(2)}|^2 \right).$$

PROOF. Let $\theta^{(i)} = \text{vec}(\otimes_{k=1}^K a_k^{(i)})$ and $\Sigma_\theta = I + \lambda^2 \theta \theta^\top$. For T i.i.d. observations \mathcal{X}_t , $t = 1, \dots, T$, the Kullback-Leibler divergence is just T times the Kullback-Leibler divergence for a single observation. Therefore, without loss of generality we take $T = 1$. Since

$$\Sigma_\theta^{-1} = I - \frac{\lambda^2}{1 + \lambda^2} \theta \theta^\top,$$

the log-likelihood function for a single observation is given by

$$\begin{aligned}\log h(x|\theta) &= -\frac{d}{2} \log(2\pi) - \frac{1}{2} \log \det(\Sigma_\theta) - \frac{1}{2} x^\top \Sigma_\theta^{-1} x \\ &= -\frac{d}{2} \log(2\pi) - \frac{1}{2} \log(1 + \lambda^2) - \frac{1}{2} x^\top x + \frac{\lambda^2}{2(1 + \lambda^2)} |x^\top \theta|^2.\end{aligned}$$

It follows that

$$\begin{aligned}D(\mathbb{P}_{\lambda, \mathbf{a}^{(1)}}^T \| \mathbb{P}_{\lambda, \mathbf{a}^{(2)}}^T) &= \mathbb{E}_{\theta^{(1)}} \left(\log h(\mathcal{X}|\theta^{(1)}) - \log h(\mathcal{X}|\theta^{(2)}) \right) \\ &= \frac{\lambda^2}{2(1 + \lambda^2)} \mathbb{E}_{\theta^{(1)}} \left(|\text{vec}(\mathcal{X})^\top \theta^{(1)}|^2 - |\text{vec}(\mathcal{X})^\top \theta^{(2)}|^2 \right) \\ &= \frac{\lambda^4}{2(1 + \lambda^2)} \left(1 - |\theta^{(1)\top} \theta^{(2)}|^2 \right) \\ &= \frac{\lambda^4}{2(1 + \lambda^2)} \left(1 - \prod_{k=1}^K |a_k^{(1)\top} a_k^{(2)}|^2 \right).\end{aligned}$$

□

Then, applying the same argument in the proof of Theorem 3.1 in [Birbaum et al. \(2013\)](#), we can obtain the desired statistical lower bound.

APPENDIX F: PROOFS OF COROLLARIES

PROPOSITION F.1. *Let $\lambda = \prod_{k=1}^K \|A_k\|_S$. Assume that the condition numbers of $A_k^\top A_k$ ($k = 1, \dots, K$) are bounded. Then, for all $1 \leq k \leq K$, we have,*

$$\begin{aligned}\|\Theta_{k,0}\|_{\text{op}} &\asymp \lambda^2 \|\Phi_{k,0}\|_{\text{op}}, & \|\Theta_{k,0}^*\|_S &\asymp \lambda^2 \|\Phi_{k,0}^*\|_S, \\ \tau_{k,r_k} &\asymp \lambda^2 \times \sigma_{\tau_k}(\text{mat}_1(\Phi_{k,1:h_0})), \\ \tau_{k,r_k}^* &\asymp \lambda^2 \times \sigma_{\tau_k}(\Phi_{k,1:h_0}^{\text{cano}}/\lambda^2).\end{aligned}$$

PROOF. If the condition numbers of $A_k^\top A_k$ ($k = 1, \dots, K$) are bounded, all the singular values of A_k are at the same order. Then Proposition F.1 immediately follows. □

PROOFS OF COROLLARY 3.1 AND 3.2. Employing Proposition F.1, under Assumption 2 and $\mathbb{E}[\text{mat}_1(\Phi_{k,1:h_0})]$ is of rank r_k , we can show $\lambda_k \asymp \lambda$ and $\tau_{k,r_k} \asymp \lambda^2$. When the ranks r_k are fixed, the second part of condition (3.10) can be written as $C_{1,K}^{(\text{iter})} R^{(\text{ideal})} \leq \rho < 1$. Thus, for $C_1^{(\text{TOPUP})} = 1/(6R^{(0)}) \leq (1 - \rho)/4$ and $C_{1,K}^{(\text{iter})} = \rho/R^{(\text{ideal})} = 1/R^{(0)}$, we have

$$\rho = C_{1,K}^{(\text{iter})} R^{(\text{ideal})} = R^{(\text{ideal})}/R^{(0)} \asymp 1/(\max_k \sqrt{d-k}).$$

For $m = 1$, this gives the rate $R^{(\text{ideal})}$ by (3.11). Then Corollary 3.1 follows from the results of Theorem 3.1.

Similarly, Applying Proposition F.1, under Assumption 2 and $\mathbb{E}[\Phi_{k,1:h_0}^{\text{cano}}/\lambda^2]$ is of rank r_k , we can obtain $\lambda_k^* \asymp \lambda$ and $\tau_{k,r_k}^* \asymp \lambda^2$. Then Corollary 3.2 follows from the results of Theorem 3.2. □

PROOFS OF COROLLARIES 3.3, 3.4 AND 3.5. By Assumption 3, (3.10) holds when

$$T \geq C_0 \max_k \left(\frac{d^{2\delta_1 - \delta_0} r^2}{r_k} + \frac{d^{2\delta_1} r r_k}{d_k} + \frac{d_k^{1/2} \sqrt{r r_k}}{d^{1/2 - \delta_1}} + \frac{d_{-k}^* \sqrt{r r_k}}{d^{1 - \delta_1}} + \frac{\sqrt{d_k d_{-k}^*} r}{d^{1 - \delta_1}} \right).$$

Because $(d_j r_j \sqrt{r r_k} + \sqrt{d_k d_j r_j} r_k)/d + \sqrt{d_k r r_k}/d r_k \leq 3r^2$ for $j \neq k$ and $\delta_1 \geq \delta_0$, the last three terms on the right-hand side above can be absorbed into the first, so that (3.25) suffices. Therefore, Theorem 3.1 yields Corollary 3.3 by setting $\rho = (R^{(\text{ideal})} + R^{(\text{add})})/R^{(0)}$ as in the proof of Corollary 3.1.

Similarly, setting $\rho = \min_k (R^{(\text{ideal})} + R^{(\text{add})}) \lambda_k^{*2} / (R^{*(0)} \|\Theta_{k,0}^*\|_S)$, and applying Assumption 3 to Theorem 3.2, leads to Corollary 3.4.

For Corollary 3.5, condition (3.20) holds when T is no smaller than

$$C_0 \max_{1 \leq k \leq K} \left(\frac{d_k r_k r_{-k}^{2\delta_2}}{d^{1 + \delta_0 - 2\delta_1}} + \frac{r_k^2 r_{-k}^{2\delta_2}}{d^{1 - 2\delta_1}} + \frac{d_k r_{-k}^2 r}{d^{1 + \delta_0 - 2\delta_1}} + \frac{d_k r_{-k}^2 r r_k}{d^{2 - 2\delta_1}} + \frac{d_{-k}^* \sqrt{r r_k}}{d^{1 - \delta_1}} + \frac{\sqrt{d_k d_{-k}^*} r}{d^{1 - \delta_1}} \right)$$

due to $d_k \sqrt{r_{-k}} \sqrt{r r_k} / d^{1 - \delta_1} \leq d_k r_k r_{-k} / d^{1 + \delta_0 - 2\delta_1}$. However, the sixth term on the right-hand side above can be absorbed into the first due to $\sqrt{d_k d_j r_j} r \leq (\sqrt{d_k} r_{-k}) (\sqrt{d_j} r_j r_{-j}) \leq \max_k d_k r_k r_{-k}^2$ for $j \neq k$. \square

APPENDIX G: TECHNICAL LEMMAS

LEMMA G.1. *Let $d, d_j, d_*, r \leq d \wedge d_j$ be positive integers, $\epsilon > 0$ and $N_{d,\epsilon} = \lfloor (1 + 2/\epsilon)^d \rfloor$. (i) For any norm $\|\cdot\|$ in \mathbb{R}^d , there exist $M_j \in \mathbb{R}^d$ with $\|M_j\| \leq 1$, $j = 1, \dots, N_{d,\epsilon}$, such that $\max_{\|M\| \leq 1} \min_{1 \leq j \leq N_{d,\epsilon}} \|M - M_j\| \leq \epsilon$. Consequently, for any linear mapping f and norm $\|\cdot\|_*$,*

$$\sup_{M \in \mathbb{R}^d, \|M\| \leq 1} \|f(M)\|_* \leq 2 \max_{1 \leq j \leq N_{d,1/2}} \|f(M_j)\|_*.$$

(ii) Given $\epsilon > 0$, there exist $U_j \in \mathbb{R}^{d \times r}$ and $V_{j'} \in \mathbb{R}^{d' \times r}$ with $\|U_j\|_S \vee \|V_{j'}\|_S \leq 1$ such that

$$\max_{M \in \mathbb{R}^{d \times d'}, \|M\|_S \leq 1, \text{rank}(M) \leq r} \min_{j \in N_{dr,\epsilon/2}, j' \in N_{d'r,\epsilon/2}} \|M - U_j V_{j'}^\top\|_S \leq \epsilon.$$

Consequently, for any linear mapping f and norm $\|\cdot\|_*$ in the range of f ,

$$(G.1) \quad \sup_{\substack{M, \tilde{M} \in \mathbb{R}^{d \times d'}, \|M - \tilde{M}\|_S \leq \epsilon \\ \|M\|_S \vee \|\tilde{M}\|_S \leq 1 \\ \text{rank}(M) \vee \text{rank}(\tilde{M}) \leq r}} \frac{\|f(M - \tilde{M})\|_*}{\epsilon^{2I_{r < d \wedge d'}}} \leq \sup_{\substack{\|M\|_S \leq 1 \\ \text{rank}(M) \leq r}} \|f(M)\|_* \leq 2 \max_{\substack{1 \leq j \leq N_{dr,1/8} \\ 1 \leq j' \leq N_{d'r,1/8}}} \|f(U_j V_{j'}^\top)\|_*.$$

(iii) Given $\epsilon > 0$, there exist $U_{j,k} \in \mathbb{R}^{d_k \times r_k}$ and $V_{j',k} \in \mathbb{R}^{d'_k \times r_k}$ with $\|U_{j,k}\|_S \vee \|V_{j',k}\|_S \leq 1$ such that

$$\max_{\substack{M_k \in \mathbb{R}^{d_k \times d'_k}, \|M_k\|_S \leq 1 \\ \text{rank}(M_k) \leq r_k, \forall k \leq K}} \min_{\substack{j_k \leq N_{d_k r_k, \epsilon/2} \\ j'_k \leq N_{d'_k r_k, \epsilon/2}, \forall k \leq K}} \left\| \bigcirc_{k=2}^K M_k - \bigcirc_{k=2}^K (U_{j_k, k} V_{j'_k, k}^\top) \right\|_{\text{op}} \leq \epsilon(K-1).$$

For any linear mapping f and norm $\|\cdot\|_*$ in the range of f ,

$$(G.2) \quad \sup_{\substack{M_k, \tilde{M}_k \in \mathbb{R}^{d_k \times d'_k}, \|M_k - \tilde{M}_k\|_S \leq \epsilon \\ \text{rank}(M_k) \vee \text{rank}(\tilde{M}_k) \leq r_k \\ \|M_k\|_S \vee \|\tilde{M}_k\|_S \leq 1 \quad \forall k \leq K}} \frac{\|f(\bigcirc_{k=2}^K M_k - \bigcirc_{k=2}^K \tilde{M}_k)\|_*}{\epsilon(2K-2)} \leq \sup_{\substack{M_k \in \mathbb{R}^{d_k \times d'_k} \\ \text{rank}(M_k) \leq r_k \\ \|M_k\|_S \leq 1, \forall k}} \|f(\bigcirc_{k=2}^K M_k)\|_*$$

and

$$(G.3) \quad \sup_{\substack{M_k \in \mathbb{R}^{d_k \times d'_k}, \|M_k\|_S \leq 1 \\ \text{rank}(M_k) \leq r_k \quad \forall k \leq K}} \|f(\bigcirc_{k=2}^K M_k)\|_* \leq 2 \max_{\substack{1 \leq j_k \leq N_{d_k r_k, 1/(8K-8)} \\ 1 \leq j'_k \leq N_{d'_k r_k, 1/(8K-8)}}} \|f(\bigcirc_{k=2}^K U_{j_k, k} V_{j'_k, k}^\top)\|_*.$$

PROOF. (i) The covering number \widetilde{N}_ϵ follows from the standard volume comparison argument as the $(1 + \epsilon/2)$ -ball under $\|\cdot\|$ and centered at the origin contains no more than $(1 + 2/\epsilon)^d$ disjoint $(\epsilon/2)$ -balls centered at M_j . The inequality follows from the ‘‘subtraction argument’’,

$$\sup_{\|M\| \leq 1} \|f(M)\|_* - \max_{1 \leq j \leq N_{d,1/2}} \|f(M_j)\|_* \leq \sup_{\|M - M_j\| \leq 1/2} \|f(M - M_j)\|_* \leq \sup_{\|M\| \leq 1} \|f(M)\|_*/2.$$

(ii) The covering numbers are given by applying (i) to both U and V in the decomposition $M = UV^\top$ as Lemma 7 in [Zhang and Xia \(2018\)](#). The first inequality in (G.1) follows from the fact that for $r < d \wedge d'$, $(M - \widetilde{M})/\epsilon$ is a sum of two rank- r matrices with no greater spectrum norm than 1, and the second inequality of (G.1) again follows from the subtraction argument although we need to split $M - U_j V_j^\top$ into two rank r matrices to result in an extra factor of 2.

(iii) The proof is nearly identical to that of part (ii). The only difference is the factor $K - 1$ when $\|\odot_{k=2}^K M_k - \odot_{k=2}^K \widetilde{M}_k\|_{\text{op}} \leq (K - 1) \max_{2 \leq k \leq K} \|M_k - \widetilde{M}_k\|_{\text{S}}$ is applied. \square

LEMMA G.2. (i) Let $G \in \mathbb{R}^{d_1 \times n}$ and $H \in \mathbb{R}^{d_2 \times n}$ be two centered independent Gaussian matrices such that $\mathbb{E}(u^\top \text{vec}(G))^2 \leq \sigma^2 \forall u \in \mathbb{R}^{d_1 n}$ and $\mathbb{E}(v^\top \text{vec}(H))^2 \leq \sigma^2 \forall v \in \mathbb{R}^{d_2 n}$. Then,

$$\|GH^\top\|_{\text{S}} \leq \sigma^2 (\sqrt{d_1 d_2} + \sqrt{d_1 n} + \sqrt{d_2 n}) + \sigma^2 x (x + 2\sqrt{n} + \sqrt{d_1} + \sqrt{d_2})$$

with at least probability $1 - 2e^{-x^2/2}$ for all $x \geq 0$.

(ii) Let $G_i \in \mathbb{R}^{d_1 \times d_2}$, $H_i \in \mathbb{R}^{d_3 \times d_4}$, $i = 1, \dots, n$, be independent centered Gaussian matrices such that $\mathbb{E}(u^\top \text{vec}(G_i))^2 \leq \sigma^2 \forall u \in \mathbb{R}^{d_1 d_2}$ and $\mathbb{E}(v^\top \text{vec}(H_i))^2 \leq \sigma^2 \forall v \in \mathbb{R}^{d_3 d_4}$. Then,

$$\begin{aligned} \left\| \text{mat}_1 \left(\sum_{i=1}^n G_i \otimes H_i \right) \right\|_{\text{S}} &\leq \sigma^2 (\sqrt{d_1 n} + \sqrt{d_1 d_3 d_4} + \sqrt{nd_2 d_3 d_4}) \\ &\quad + \sigma^2 x (x + \sqrt{n} + \sqrt{d_1} + \sqrt{d_2} + \sqrt{d_3 d_4}) \end{aligned}$$

with at least probability $1 - 2e^{-x^2/2}$ for all $x \geq 0$.

PROOF. Assume $\sigma = 1$ without loss of generality. Let $x \geq 0$.

(i) Independent of G and H , let $\zeta_j \in \mathbb{R}^n$, $j = 1, 2$, be independent standard Gaussian vectors. As in [Chen, Yang and Zhang \(2019\)](#), the Sudakov-Fernique inequality provides

$$\mathbb{E} \left[\|GH^\top\|_{\text{S}} \middle| G \right] \leq \mathbb{E} \left[\max_{\|u\|_2=1} u^\top G \zeta_2 \middle| G \right] + \|G\|_{\text{S}} \sqrt{d_2}.$$

Thus, by the Gaussian concentration inequality

$$\mathbb{P} \left\{ \|GH^\top\|_{\text{S}} \geq \mathbb{E} \left[\max_{\|u\|_2=1} u^\top G \zeta_2 \middle| G \right] + \|G\|_{\text{S}} (\sqrt{d_2} + x) \middle| G \right\} \leq e^{-x^2/2}.$$

Applying the Sudakov-Fernique inequality again, we have

$$\mathbb{E} \left[\mathbb{E} \left[\max_{\|u\|_2=1} u^\top G \zeta_2 \middle| G \right] + \|G\|_{\text{S}} (\sqrt{d_2} + x) \right] \leq \sqrt{d_1 n} + (\sqrt{d_1} + \sqrt{n}) (\sqrt{d_2} + x).$$

Moreover, as the Lipschitz norm of $\mathbb{E} \left[\max_{\|u\|_2=1} u^\top G \zeta_2 \middle| G \right] + \|G\|_{\text{S}} (\sqrt{d_2} + x)$ is bounded by $\sqrt{n} + \sqrt{d_2} + x$, by the Gaussian concentration inequality

$$\begin{aligned} &\mathbb{E} \left[\max_{\|u\|_2=1} u^\top G \zeta_2 \middle| G \right] + \|G\|_{\text{S}} (\sqrt{d_2} + x) \\ &\leq \sqrt{d_1 n} + (\sqrt{d_1} + \sqrt{n}) (\sqrt{d_2} + x) + x (\sqrt{n} + \sqrt{d_2} + x) \end{aligned}$$

holds with at least probability $1 - e^{-x^2/2}$.

(ii) We treat $G = (G_1, \dots, G_n) \in \mathbb{R}^{d_1 \times d_2 \times n}$ and $H = (H_1, \dots, H_n) \in \mathbb{R}^{d_3 \times d_4 \times n}$ as tensors. Let $\xi = (\xi_1, \dots, \xi_n) \in \mathbb{R}^{d_2 \times n}$ be a standard Gaussian matrix independent of H . For $u \in \mathbb{R}^{d_1}$ and $V \in \mathbb{R}^{d_2 \times (d_3 d_4)}$,

$$\begin{aligned}
& \mathbb{E} \left[\left\| \text{mat}_1 \left(\sum_{i=1}^n G_i \otimes H_i \right) \right\|_{\text{S}} \middle| H \right] \\
&= \mathbb{E} \left[\sup_{\|u\|_2=1, \|V\|_{\text{F}}=1} u^\top \text{mat}_1(G) \text{vec}(\text{mat}_3(H) V^\top) \middle| H \right] \\
&\leq \sqrt{d_1} \sup_{\|V\|_{\text{F}}=1} \|\text{mat}_3(H) V^\top\|_{\text{F}} + \mathbb{E} \left[\sup_{\|V\|_{\text{F}}=1} (\text{vec}(\xi))^\top \text{vec}(\text{mat}_3(H) V^\top) \middle| H \right] \\
&= \sqrt{d_1} \|\text{mat}_3(H)\|_{\text{S}} + \mathbb{E} \left[\left(\sum_{j=1}^{d_2} \sum_{k=1}^{d_3 d_4} \left(\sum_{i=1}^n \xi_{i,j} \text{vec}(H_i)_k \right)^2 \right)^{1/2} \middle| H \right] \\
&\leq \sqrt{d_1} \|\text{mat}_3(H)\|_{\text{S}} + \sqrt{d_2} \|\text{vec}(H)\|_2
\end{aligned}$$

By the Gaussian concentration inequality,

$$\mathbb{P} \left\{ \left\| \text{mat}_1 \left(\sum_{i=1}^n G_i \otimes H_i \right) \right\|_{\text{S}} \geq (\sqrt{d_1} + x) \|\text{mat}_3(H)\|_{\text{S}} + \sqrt{d_2} \|\text{vec}(H)\|_2 \middle| H \right\} \leq e^{-x^2/2}.$$

Moreover, as $\mathbb{E}[(\sqrt{d_1} + x) \|\text{mat}_3(H)\|_{\text{S}} + \sqrt{d_2} \|\text{vec}(H)\|_2] \leq (\sqrt{d_1} + x)(\sqrt{n} + \sqrt{d_3 d_4}) + \sqrt{d_2 n d_3 d_4}$ and the Lipschitz norm of $(\sqrt{d_1} + x) \|\text{mat}_3(H)\|_{\text{S}} + \sqrt{d_2} \|\text{vec}(H)\|_2$ is bounded by $\sqrt{d_1} + x + \sqrt{d_2}$,

$$\begin{aligned}
& (\sqrt{d_1} + x) \|\text{mat}_3(H)\|_{\text{S}} + \sqrt{d_2} \|\text{vec}(H)\|_2 \\
&\leq (\sqrt{d_1} + x)(\sqrt{n} + \sqrt{d_3 d_4}) + \sqrt{nd_2 d_3 d_4} + x(\sqrt{d_1} + x + \sqrt{d_2})
\end{aligned}$$

holds with at least probability $1 - e^{-x^2/2}$. □