# Tensor SVD: Statistical and Computational Limits \*

Anru Zhang and Dong Xia University of Wisconsin-Madison

#### Abstract

In this paper, we propose a general framework for tensor singular value decomposition (tensor SVD), which focuses on the methodology and theory for extracting the hidden low-rank structure from high-dimensional tensor data. Comprehensive results are developed on both the statistical and computational limits for tensor SVD. This problem exhibits three different phases according to the signal-to-noise ratio (SNR). In particular, with strong SNR, we show that the classical higher-order orthogonal iteration achieves the minimax optimal rate of convergence in estimation; with weak SNR, the information-theoretical lower bound implies that it is impossible to have consistent estimation in general; with moderate SNR, we show that the non-convex maximum likelihood estimation provides optimal solution, but with NP-hard computational cost; moreover, under the hardness hypothesis of hypergraphic planted clique detection, there are no polynomial-time algorithms performing consistently in general.

## 1 Introduction

There is no need to argue the importance of singular value decomposition (SVD) in data analysis. As one of the most important tools in multivariate analysis, SVD along with the closely related formulation, i.e., principal component analysis (PCA), have been a mainstay of data analysis

<sup>\*</sup>Anru Zhang is Assistant Professor, Department of Statistics, University of Wisconsin-Madison, Madison, WI 53706, E-mail: anruzhang@stat.wisc.edu; Dong Xia is Visiting Assistant Professor, Department of Statistics, University of Wisconsin-Madison, Madison, WI 53706, E-mail: dongxia@stat.wisc.edu.

since more than a century ago, and widely used in various subjects. Attributed to the modern high-dimensional data, the popularity of SVD and PCA continues to surge in the recent decades, and many important variations, such as sparse SVD [1, 2, 3, 4], matrix denoising [5, 6, 7, 8], sparse PCA [9, 10, 11], robust PCA [12], have been proposed and developed recently. Traditionally, most of the SVD and PCA results focused on exploiting low-rank structures from datasets in form of matrices.

Motivated by modern scientific research, tensors, or high-order arrays, have been actively studied in machine learning, electrical engineering, and statistics. Some specific scientific applications involving tensor data include neuroimaging analysis [13, 14], recommender systems [15, 16], computer vision [17, 18], topic modeling [19], community detection [20], hyperspectral image compression [21], spatiotemporal gene expression [22], etc. A common objective in these problems is to dig out the underlying high-order low-rank structure, such as the singular subspaces and the whole low-rank tensors, buried in the noisy observations. To achieve this goal, we are in strong need of a statistical tool for tensor data that is the counterpart of regular singular value decomposition for traditional order-2 datasets. Richard and Montanari [23], Hopkins et al [24], Perry et al [25] considered a rank-1 spiked tensor SVD statistical model and proposed various methods, including tensor unfolding and sum of square optimization (SOS). However, as far as we know, the statistical framework for general rank-r high-order tensor SVD or PCA was not well established or studied in the literature.

In this paper, we propose a general framework of tensor singular value decomposition (tensor SVD). To be specific, suppose we are interested in a low-rank tensor  $\mathbf{X} \in \mathbb{R}^{p_1 \times p_2 \times p_3}$ , which is observed with entry-wise corruptions as follows,

$$Y = X + Z. (1)$$

Here **Z** is the  $p_1$ -by- $p_2$ -by- $p_3$  noisy tensor with  $\{Z_{ijk}\}_{i,j,k=1}^{p_1,p_2,p_3} \stackrel{iid}{\sim} N(0,\sigma^2)$ ; **X** is a fixed tensor with low Tucker ranks in the sense that all fibers of **X** along three directions (i.e., counterpart of matrix columns and rows for tensors, also see Section 2.1 for formal definitions) lie in low-dimensional subspaces, say  $U_1$ ,  $U_2$ , and  $U_3$ , respectively. Our goal is to estimate  $U_1$ ,  $U_2$ ,  $U_3$ , and **X** from the noisy observation **Y**.

It is worth mentioning that the analog of this problem when X is an order-2 tensor, i.e.,

a matrix, has been previously studied in the context of matrix denoising in [5, 7, 8, 26]. For the matrix denoising problem, the best low-rank matrix approximation provides the optimal results, which can be calculated efficiently via singular value decomposition, as guaranteed by the well-regarded Eckart-Young-Mirsky Theorem.

Although there have been significant efforts in developing methodologies and theories for matrix SVD or matrix denoising, there is a paucity of literature on the analogous question for tensors of order three or higher. In fact, SVD for high-order tensors is much more difficult than its counterpart for matrices in various of aspects. First, tensors have more involved structures along three or more ways, while the traditional tools for matrices could typically incorporate two ways. As we will see later, one may achieve a sub-optimal result by simply ignoring the structure beyond two ways. Second, many operations for matrices, such as operator norm, singular value decomposition, are either not well defined or computational NP-hard for high-order tensors [27]. Third, high-order tensors often bring about high-dimensionality. For incidence, a 500-by-500-by-500 tensor is comprised of more than 12,500,000 entries that impose significant computational challenges. All these characteristics make the tensor SVD distinct from the classical matrix problems.

The best low-rank tensor approximation, or equivalently the maximum likelihood estimation (MLE), is a straightforward solution for tensor SVD. However, MLE is non-convex and computationally NP-hard in general (see, e.g. Hillar and Lim [27]). De Lathauwer, De Moor, and Vandewalle instead introduced the higher order SVD (HOSVD) [28] and higher order orthogonal iteration (HOOI) [29], which aims at approximating the best low-rank approximation with efficient spectral and power iteration method. Since then, HOSVD and HOOI have been widely studied in the literature (see, e.g. [30, 31, 32, 33, 34]). However as far as we know, many basic theoretical properties of these procedures, such as the error bound and the necessary iteration times, still remain unclear.

In this paper, we develop comprehensive results on both the statistical and computational limits for tensor SVD. To be specific, we establish upper bounds on estimation errors for both higher-order orthogonal iteration (HOOI) and maximum likelihood estimator (MLE). It is also shown that HOOI converges within a logarithm factor of iterations. Then the matching

information-theoretical lower bounds over a large class of low-rank tensors are correspondingly introduced. To the best of our knowledge, we are among the first to develop the statistical guarantees for both HOOI and MLE. Let the Tucker rank of  $\mathbf{X}$  be  $(r_1, r_2, r_3)$  (see formal definition in Section 3). The statistical and computational barriers of tensor SVD problem rely on a key factor  $\lambda$ , i.e., the smallest non-zero singular values of matricizations of  $\mathbf{X}$  (also see formal definition in Section 3), which essentially measures the signal strength of the problem. When  $p = \min\{p_1, p_2, p_3\}$ ,  $p_k \leq Cp$ ,  $r_k \leq Cp^{1/2}$  for k = 1, 2, 3 and a constant C > 0, our main results can be summarized into the following three phases according to signal-to-noise ratio (SNR):  $\lambda/\sigma$ .

- 1. When  $\lambda/\sigma = p^{\alpha}$  for  $\alpha \geq 3/4$ , the scenario is referred to as the **strong SNR case**. The fast higher-order orthogonal iteration (HOOI) recovers  $U_1, U_2, U_3$ , and **X** with the minimax optimal rate of convergence over a general class of low-rank tensors.
- 2. When  $\lambda/\sigma = p^{\alpha}$  for  $\alpha < 1/2$ , we refer to this case as the **weak SNR case**, and propose the minimax lower bound to show that there are no consistent estimators of  $U_1, U_2, U_3$ , or  $\mathbf{X}$ ;
- 3. When  $\lambda/\sigma = p^{\alpha}$  for  $1/2 \leq \alpha < 3/4$ , the scenario is referred to as the **moderate SNR** case. We provide a computational lower bound to show that no polynomial time algorithm can recover  $U_1, U_2, U_3$  consistently based on an assumption of hypergraphic planted clique detection. Meanwhile, the maximum likelihood estimator, although being computational intractable, achieves optimal rates of convergence over a general class of low-rank tensors.

It is also noteworthy that our results can be further generalized to fourth or higher order tensors, or when the noise **Z** is i.i.d. sub-Gaussian distributed.

Our work is also related to several recent results in literature. For example, [23, 24, 35, 36] considered the extraction of rank-1 symmetric tensors from i.i.d. (symmetric) Gaussian noise, which is a rank-1 special case of our tensor SVD model; [37, 38] considered the CP low-rank tensor decomposition based on noisy observations; [25] considered the statistical limit of detecting and estimating a randomly sampled rank-one structure from a symmetric random Gaussian tensor; [39, 40] considered the regularized tensor factorizations with/without sparsity; [41] and [42] further considered non-negative tensor decomposition and robust tensor principal

component analysis; [22] focused on orthogonal decomposable tensor SVD problem; Lesieur et al [43] considered a Bayesian symmetric spiked tensor estimation model – an approximate message passing algorithm (AMP) was particularly introduced and the rigorous asymptotic analysis for statistical and computational phase transitions were performed on high-order, symmetric, and rank-1 tensor estimation. It should be noted that different from previous works, we perform non-asymptotic analysis for tensor SVD, where the signal tensor  $\mathbf{X}$  can be generally Tuckerrank-r, non-random, and asymmetric. Also, to the best of our knowledge, we are among the first to provide a comprehensive analysis of both statistical and computational optimality of tensor SVD.

The rest of the article is organized as follows. After a brief explanation for basic notations and tensor algebra in Section 2.1, we state the fast higher-order orthogonal iteration and the non-convex maximum likelihood estimation for tensor SVD in Section 2.2. The statistical limits in the context of minimax optimality are provided for strong, weak, and moderate SNR cases respectively in Section 3. Then we further discuss the computational barriers in moderate SNR case in Section 4. Simulation results are provided in Section 5 to justify the theoretical results of this paper. We briefly discuss the extension of the results to fourth or higher order tensors and i.i.d. sub-Gaussian noise cases in Section 6. The proofs of all technical results are given in Section 7 and the supplementary materials.

## 2 Tensor SVD: Methodology

#### 2.1 Notations, Preliminaries, and Tensor Algebra

In this section, we start with basic notations and tensor algebra to be used throughout the paper. For  $a,b \in \mathbb{R}$ , let  $a \wedge b = \min\{a,b\}$ ,  $a \vee b = \max\{a,b\}$ . For two sequences  $\{a_i\},\{b_i\}$ , if there are two constants C,c>0 such that  $ca_i \leq b_i \leq Ca_i$  for all  $i \geq 1$ , we denote  $a \times b$ .  $C,c,C_0,c_0,\ldots$  represent generic constants, whose actual values of these generic constants may vary from time to time. Particularly, the uppercase and lowercase letters represent large and small constants, respectively. The matrices are denoted as capital letters,  $U_1,V_1,A$ , etc. Especially,  $\mathbb{O}_{p,r} := \{U \in \mathbb{R}^{p \times r} : U^{\top}U = I_r\}$  is the set of all p-by-r matrices with or-

thonormal columns. For any matrix  $A \in \mathbb{R}^{p_1 \times p_2}$ , let  $\sigma_1(A) \geq \cdots \geq \sigma_{p_1 \wedge p_2}(A) \geq 0$  be the singular values in non-increasing order. We are particularly interested in the smallest singular value of A:  $\sigma_{\min}(A) = \sigma_{p_1 \wedge p_2}(A)$ . In addition, the class of matrix Schatten q-norms will be used:  $\|A\|_q = \left(\sum_{j=1}^{p_1 \wedge p_2} \sigma_j^q(A)\right)^{1/q}$ . Specific instances of Schatten q-norms include the Frobenius norm (i.e., Schatten 2-norm),  $\|A\|_F = \sqrt{\sum_{i=1}^{p_1} \sum_{j=1}^{p_2} A_{ij}^2} = \sqrt{\sum_{j=1}^{p_1 \wedge p_2} \sigma_j^2(A)}$ , and spectral norm (i.e., Schatten  $\infty$ -norm),  $\|A\| = \sigma_1(A) = \max_{v \in \mathbb{R}^{p_2}} \frac{\|Av\|_2}{\|v\|_2}$ . We also use  $\mathrm{SVD}_r(A)$  to denote the leading r left singular vectors of A, so that  $\mathrm{SVD}_r(A) \in \mathbb{O}_{p_1,r}$ . Define the projection operator  $P_A = A(A^\top A)^\dagger A^\top$ . Here  $(\cdot)^\dagger$  represents the psudo-inverse. If  $A = U\Sigma V^\top$  is the  $\mathrm{SVD}$ ,  $P_A$  can be equivalently written as  $P_A = UU^\top$ . For any two matrices, say  $U \in \mathbb{R}^{p_1 \times r_1}$ ,  $V \in \mathbb{R}^{p_2 \times r_2}$ , we also let  $U \otimes V \in \mathbb{R}^{(p_1 p_2) \times (r_1 r_2)}$  be their outer product matrix, such that  $(U \otimes V)_{[(i-1)p_3+j,(k-1)r_3+l]} = U_{ik} \cdot V_{jl}$ , for  $i=1,\ldots,p_1,j=1,\ldots,p_2,k=1,\ldots,r_1$ , and  $l=1,\ldots,r_2$ . We adopt R convention to represent submatrices:  $A_{[a:b,c:d]}$  represents the a-to-b-th rows, c-to-d-th columns of matrix A; we also use  $A_{[a:b,c:d]}$  to represent a-to-b-th full rows of A and C-to-d-th full columns of A, respectively.

We use  $\sin \Theta$  distances to measure the difference between singular subspaces. To be specific, for any two  $p \times r$  matrices with orthonormal columns, say U and  $\hat{U}$ , we define the principal angles between U and  $\hat{U}$  as  $\Theta(U, \hat{U}) = \operatorname{diag}(\arccos(\sigma_1), \dots, \arccos(\sigma_r)) \in \mathbb{R}^{r \times r}$ , where  $\sigma_1 \geq \dots \geq \sigma_r \geq 0$  are the singular values of  $U^{\top}\hat{U}$ . The Schatten q-sin  $\Theta$ -norm is then defined as

$$\|\sin\Theta(U,\hat{U})\|_q = \left(\sum_{i=1}^r \sin^q(\arccos(\sigma_i))\right)^{1/q} = \left(\sum_{i=1}^r \left(1 - \sigma_i^2\right)^{q/2}\right)^{1/q}, \quad 1 \le q \le +\infty.$$

The readers are referred to Lemma 3 in the supplementary materials and Lemma 1 in [26] for more discussions on basic properties of  $\sin \Theta$  distances.

Throughout this paper, the boldface capital letters, e.g.  $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$ , note tensors. To simplify the presentation, the main context of this paper is focused on third order tensor. The extension to 4-th or higher tensors is briefly discussed in Section 6. The readers are also referred to [44] for a more detailed tutorial of tensor algebra. For any tensor  $\mathbf{X} \in \mathbb{R}^{p_1 \times p_2 \times p_3}$ , define its mode-1 matricization as a  $p_1$ -by- $(p_2p_3)$  matrix  $\mathcal{M}_1(\mathbf{X})$  such that

$$[\mathcal{M}_1(\mathbf{X})]_{i,(j-1)p_3+k} = X_{ijk}, \quad \forall 1 \le i \le p_1, 1 \le j \le p_2, 1 \le k \le p_3.$$

In other words,  $\mathcal{M}_1(\mathbf{X})$  is composed of all mode-1 fibers,  $\{(\mathbf{X}_{[:,i_2,i_3]}) \in \mathbb{R}^{p_1} : 1 \leq i_2 \leq p_2, 1 \leq i_3 \leq p_3\}$ , of  $\mathbf{X}$ . The mode-2 and mode-3 matricizations, i.e.,  $\mathcal{M}_2(\mathbf{X}) \in \mathbb{R}^{p_2 \times (p_3 p_1)}$  and  $\mathcal{M}_3 \in \mathcal{M}_2(\mathbf{X})$ 

 $\mathbb{R}^{p_3 \times (p_1 p_2)}$ , are defined in the same fashion. We also define the marginal multiplication  $\times_1$ :  $\mathbb{R}^{p_1 \times p_2 \times p_3} \times \mathbb{R}^{r_1 \times p_1} \to \mathbb{R}^{r_1 \times p_2 \times p_3}$  as

$$\mathbf{X} \times_1 Y = \left( \sum_{i'=1}^{p_1} X_{i'jk} Y_{i,i'} \right)_{1 \le i \le r_1, 1 \le j \le p_2, 1 \le k \le p_3}.$$

Marginal multiplications  $\times_2$  and  $\times_3$  can be defined similarly.

Different from matrices, there is no universal definition for tensor ranks. We particularly introduce the following Tucker ranks (also called multilinear ranks) of X as

$$r_1 = \operatorname{rank}_1(\mathbf{X}) = \operatorname{rank}(\mathcal{M}_1(\mathbf{X}))$$
  
=  $\dim(\operatorname{span}\{\mathbf{X}_{[:,i_2,i_3]} \in \mathbb{R}^{p_1} : 1 \le i_2 \le p_2, 1 \le i_3 \le p_3\}).$ 

 $r_2 = \operatorname{rank}_2(\mathbf{X})$  and  $r_3 = \operatorname{rank}_3(\mathbf{X})$  can be similarly defined. Note that, in general,  $r_1, r_2, r_3$  satisfy  $r_1 \leq r_2 r_3, r_2 \leq r_3 r_1, r_3 \leq r_1 r_2$ , but are not necessarily equal. We further denote  $\operatorname{rank}(\mathbf{X})$  as the triplet:  $(r_1, r_2, r_3)$ . The Tucker  $\operatorname{rank}(r_1, r_2, r_3)$  is also closely associated with the following Tucker decomposition. Let  $U_1 \in \mathbb{O}_{p_1, r_1}, U_2 \in \mathbb{O}_{p_2, r_2}, U_3 \in \mathbb{O}_{p_3, r_3}$  be the left singular vectors of  $\mathcal{M}_1(\mathbf{X})$ ,  $\mathcal{M}_2(\mathbf{X})$  and  $\mathcal{M}_3(\mathbf{X})$  respectively, then there exists a core tensor  $\mathbf{S} \in \mathbb{R}^{r_1 \times r_2 \times r_3}$  such that

$$\mathbf{X} = \mathbf{S} \times_1 U_1 \times_2 U_2 \times_3 U_3, \quad \text{or} \quad X_{ijk} = \sum_{i'=1}^{r_1} \sum_{j'=1}^{r_2} \sum_{k'=1}^{r_3} S_{i'j'k'}(U_1)_{i,i'}(U_2)_{j,j'}(U_3)_{k,k'}. \tag{2}$$

Expression (2) is widely referred to as the Tucker decomposition of  $\mathbf{X}$ . Finally, to measure the tensor estimation error, we introduce the following tensor Frobenius norm,

$$\|\mathbf{X}\|_{\mathrm{F}} = \Big(\sum_{i=1}^{p_1} \sum_{j=1}^{p_2} \sum_{k=1}^{p_3} X_{ijk}^2\Big)^{1/2}.$$

### 2.2 Maximum Likelihood Estimator and Higher-Order Orthogonal Iteration

In this section, we discuss the methodology for tensor SVD. Given the knowledge of Tucker decomposition, the original tensor SVD model (1) can be cast as follows,

$$\mathbf{Y} = \mathbf{X} + \mathbf{Z} = \mathbf{S} \times_1 U_1 \times_2 U_2 \times_3 U_3 + \mathbf{Z}, \quad \mathbf{Z} \stackrel{i.i.d.}{\sim} N(0, \sigma^2), \tag{3}$$

where  $U_1 \in \mathbb{O}_{p_1,r_1}$ ,  $U_2 \in \mathbb{O}_{p_2,r_2}$ ,  $U_3 \in \mathbb{O}_{p_3,r_3}$ , and  $\mathbf{S} \in \mathbb{R}^{r_1 \times r_2 \times r_3}$ . Our goal is to estimate  $U_1, U_2, U_3$ , and  $\mathbf{X}$  from  $\mathbf{Y}$ . Clearly, the log-likelihood of Model (1) can be written (ignoring

the constants) as  $\mathcal{L}(\mathbf{Y}|\mathbf{X}) = -\frac{1}{\sigma^2} \|\mathbf{Y} - \mathbf{X}\|_F^2$ , then it is straightforward to apply the maximum likelihood estimator for estimation,

$$\hat{\mathbf{X}}^{\text{mle}} = \underset{\text{rank}(\mathbf{X}) \leq r_1, r_2, r_3}{\arg \min} \|\mathbf{Y} - \mathbf{X}\|_{\text{F}},$$

$$\hat{U}_k^{\text{mle}} = \text{SVD}_{r_k}(\hat{\mathbf{X}}^{\text{mle}}), \quad k = 1, 2, 3.$$

Intuitively speaking, MLE seeks the best rank- $(r_1, r_2, r_3)$  approximation for **Y** in Frobenius norm. By Theorems 4.1 and 4.2 in [29], MLE can be equivalently written as

$$\hat{U}_{1}^{\text{mle}}, \hat{U}_{2}^{\text{mle}}, \hat{U}_{3}^{\text{mle}} = \underset{V_{k} \in \mathbb{O}_{p_{k}, r_{k}}}{\operatorname{arg\,max}} \left\| \mathbf{Y} \times_{1} V_{1}^{\top} \times_{2} V_{2}^{\top} \times_{3} V_{3}^{\top} \right\|_{F}^{2},$$

$$\hat{\mathbf{X}}^{\text{mle}} = \mathbf{Y} \times_{1} P_{\hat{U}_{2}^{\text{mle}}} \times_{2} P_{\hat{U}_{2}^{\text{mle}}} \times_{3} P_{\hat{U}_{2}^{\text{mle}}}.$$
(4)

As we will illustrate later in Section 3, such estimators achieve optimal rate of convergence in estimation errors. On the other hand, (4) is non-convex and computationally NP-hard even when r = 1 (see, e.g., [27]). Then MLE may not be applicable in practice.

To overcome the computational difficulties of MLE, we consider a version of higher-order orthogonal iteration (HOOI) [29]. The procedure includes three steps: spectral initialization, power iteration and tensor projection. The first two steps produce optimal estimation of loadings  $U_1, U_2, U_3$ . The final step yields an optimal estimator of the underlying low rank tensor **X**. It is helpful for us to present the procedure of HOOI in details as follows.

Step 1 (Spectral initialization) Since  $U_1, U_2$ , and  $U_3$  respectively represent the singular subspaces of  $\mathcal{M}_1(\mathbf{X})$ ,  $\mathcal{M}_2(\mathbf{X})$ , and  $\mathcal{M}_3(\mathbf{X})$ , it is natural to perform singular value decomposition (SVD) on  $\mathcal{M}_k(\mathbf{Y})$  to obtain preliminary estimators for  $U_k$ :

$$\hat{U}_k^{(0)} = \text{SVD}_{r_k}(\mathcal{M}_k(\mathbf{Y})) = \text{the first } r_k \text{ left singular vectors of } \mathcal{M}_k(\mathbf{Y}).$$

In fact,  $\hat{U}_k^{(0)}$  is exactly the higher-order SVD (HOSVD) estimator introduced by De Lathauwer, De Moor, and Vandewalle [28]. As we will show later,  $\hat{U}_k^{(0)}$  serves as a good starting point but not an optimal estimator for  $U_k$ .

Step 2 (Power Iteration) Then one applies power iterations to update the estimations. Given  $\hat{U}_2^{(t-1)}, \hat{U}_3^{(t-1)}, \hat{V}_3$  can be denoised via mode-2 and 3 projections:  $\mathbf{Y} \times_2 (\hat{U}_2^{(t-1)})^{\top} \times_3 (\hat{U}_3^{(t-1)})^{\top}$ . As we will illustrate via theoretical analysis, the mode-1 singular subspace

of **X** is preserved while the amplitude of the noise is highly reduced after such the projection. Thus, for t = 1, 2, ..., we calculate

$$\hat{U}_{1}^{(t)} = \text{first } r_{1} \text{ left singular vectors of } \mathcal{M}_{1}(\mathbf{Y} \times_{2} (\hat{U}_{2}^{(t-1)})^{\top} \times_{3} (\hat{U}_{3}^{(t-1)})^{\top}), 
\hat{U}_{2}^{(t)} = \text{first } r_{2} \text{ left singular vectors of } \mathcal{M}_{2}(\mathbf{Y} \times_{1} (\hat{U}_{1}^{(t)})^{\top} \times_{3} (\hat{U}_{3}^{(t-1)})^{\top}), 
\hat{U}_{3}^{(t)} = \text{first } r_{3} \text{ left singular vectors of } \mathcal{M}_{3}(\mathbf{Y} \times_{1} (\hat{U}_{1}^{(t)})^{\top} \times_{2} (\hat{U}_{2}^{(t)})^{\top}).$$
(5)

The iteration is stopped when either the increment is no more than the tolerance  $\varepsilon$ , i.e.,

$$\left\| \mathbf{Y} \times_{1} (\hat{U}_{1}^{(t)})^{\top} \times_{2} (\hat{U}_{2}^{(t)})^{\top} \times_{3} (\hat{U}_{3}^{(t)})^{\top} \right\|_{F} 
- \left\| \mathbf{Y} \times_{1} (\hat{U}_{1}^{(t-1)})^{\top} \times_{2} (\hat{U}_{2}^{(t-1)})^{\top} \times_{3} (\hat{U}_{3}^{(t-1)})^{\top} \right\|_{F} \le \varepsilon,$$
(6)

or the maximum number of iterations is reached.

Step 3 (Projection) With the final estimates  $\hat{U}_1, \hat{U}_2, \hat{U}_3$ , it is natural to estimate **S** and **X** as

$$\hat{\mathbf{S}} = \mathbf{Y} \times_1 \hat{U}_1^\top \times_2 \hat{U}_2^\top \times_3 \hat{U}_3^\top, \quad \hat{\mathbf{X}} = \hat{\mathbf{S}} \times_1 \hat{U}_1 \times_2 \hat{U}_2 \times_3 \hat{U}_3 = \mathbf{Y} \times_1 P_{\hat{U}_1} \times_2 P_{\hat{U}_2} \times_3 P_{\hat{U}_3}.$$

The procedure of HOOI is summarized in Algorithm 1. The further generalization to order-4 or higher tensors SVD will be discussed in Section 6.

## 3 Statistical Limits: Minimax Upper and Lower Bounds

In this section, we develop the statistical limits for tensor SVD. Particularly, we analyze the estimation error upper bounds of HOOI and MLE, then develop the corresponding lower bounds. For any  $\mathbf{X} \in \mathbb{R}^{p_1 \times p_2 \times p_3}$ , denote  $\lambda = \min_{k=1,2,3} \sigma_{r_k}(\mathcal{M}_k(\mathbf{X}))$  as the minimal singular values of each matricization, which essentially measures the signal level in tensor SVD model. Suppose the signal-to-noise ratio (SNR) is  $\lambda/\sigma = p^{\alpha}$ , where  $p = \min\{p_1, p_2, p_3\}$ . Then the problem of tensor SVD exhibits three distinct phases:  $\alpha \geq 3/4$  (strong SNR),  $\alpha < 1/2$  (weak SNR), and  $1/2 \leq \alpha < 3/4$  (moderate SNR).

We first analyze the statistical performance of HOOI, i.e., Algorithm 1, under the strong SNR setting that  $\lambda/\sigma \geq Cp^{3/4}$ .

**Theorem 1** (Upper Bound for HOOI). Suppose there exist constants  $C_0, c_0 > 0$  such that  $p_k \leq C_0 p$ ,  $\|\mathbf{X}\|_{\mathrm{F}} \leq C_0 \sigma \exp(c_0 p)$ ,  $r_k \leq C_0 p^{1/2}$ , for  $p = \min\{p_1, p_2, p_3\}$  and k = 1, 2, 3. Then

### Algorithm 1 Higher Order Orthogonal Iteration (HOOI) [29]

- 1: Input:  $\mathbf{Y} \in \mathbb{R}^{p_1 \times p_2 \times p_3}$ ,  $(r_1, r_2, r_3)$ , increment tolerance  $\varepsilon > 0$ , and maximum number of iterations  $t_{\text{max}}$ .
- 2: Let t = 0, initiate via matricization SVDs

$$\hat{U}_{1}^{(0)} = \text{SVD}_{r_{1}}(\mathcal{M}_{1}(\mathbf{Y})), \quad \hat{U}_{2}^{(0)} = \text{SVD}_{r_{2}}(\mathcal{M}_{2}(\mathbf{Y})), \quad \hat{U}_{3}^{(0)} = \text{SVD}_{r_{3}}(\mathcal{M}_{3}(\mathbf{Y})).$$

- 3: repeat
- 4: Let t = t + 1, calculate

$$\hat{U}_{1}^{(t)} = \text{SVD}_{r_{1}} \left( \mathcal{M}_{r_{1}} (\mathbf{Y} \times_{2} (\hat{U}_{2}^{(t-1)})^{\top} \times_{3} (\hat{U}_{3}^{(t-1)})^{\top}) \right),$$

$$\hat{U}_{2}^{(t)} = \text{SVD}_{r_{2}} \left( \mathcal{M}_{r_{2}} (\mathbf{Y} \times_{1} (\hat{U}_{1}^{(t)})^{\top} \times_{3} (\hat{U}_{3}^{(t-1)})^{\top}) \right),$$

$$\hat{U}_{3}^{(t)} = \text{SVD}_{r_{3}} \left( \mathcal{M}_{r_{3}} (\mathbf{Y} \times_{1} (\hat{U}_{1}^{(t)})^{\top} \times_{2} (\hat{U}_{2}^{(t)})^{\top}) \right).$$

5: **until**  $t = t_{\text{max}}$  or

$$\begin{aligned} & \left\| \mathbf{Y} \times_{1} (\hat{U}_{1}^{(t)})^{\top} \times_{2} (\hat{U}_{2}^{(t)})^{\top} \times_{3} (\hat{U}_{3}^{(t)})^{\top} \right\|_{F} \\ & - \left\| \mathbf{Y} \times_{1} (\hat{U}_{1}^{(t-1)})^{\top} \times_{2} (\hat{U}_{2}^{(t-1)})^{\top} \times_{3} (\hat{U}_{3}^{(t-1)})^{\top} \right\|_{F} \leq \varepsilon. \end{aligned}$$

6: Estimate and output:

$$\hat{U}_1 = \hat{U}_1^{(t_{\text{max}})}, \quad \hat{U}_2 = \hat{U}_2^{(t_{\text{max}})}, \quad \hat{U}_3 = \hat{U}_3^{(t_{\text{max}})};$$

$$\hat{\mathbf{X}} = \mathbf{Y} \times_1 P_{\hat{U}_1} \times_2 P_{\hat{U}_2} \times_3 P_{\hat{U}_3}.$$

there exist absolute constants  $C_{gap}$ , C > 0, which does not depend on  $p_k$ ,  $r_k$ ,  $\lambda$ ,  $\sigma$ , q, such that whenever

$$\lambda/\sigma \ge C_{gap}p^{3/4}$$
, (i.e., in strong SNR case),

after at most  $t_{\max} = C\left(\log\left(\frac{p}{\lambda}\right) \vee 1\right)$  iterations in Algorithm 1, the following upper bounds hold,

$$\mathbb{E}r_k^{-1/q} \left\| \sin \Theta \left( \hat{U}_k, U_k \right) \right\|_q \le C \frac{\sqrt{p_k}}{\lambda/\sigma}, \quad k = 1, 2, 3, \quad 1 \le q \le \infty, \tag{7}$$

$$\mathbb{E} \| \hat{\mathbf{X}} - \mathbf{X} \|_{\mathrm{F}}^{2} \le C \sigma^{2} \left( p_{1} r_{1} + p_{2} r_{2} + p_{3} r_{3} \right), \quad \mathbb{E} \frac{\| \hat{\mathbf{X}} - \mathbf{X} \|_{\mathrm{F}}^{2}}{\| \mathbf{X} \|_{\mathrm{F}}^{2}} \le C \left( \frac{(p_{1} + p_{2} + p_{3})}{\lambda^{2} / \sigma^{2}} \bigwedge 1 \right). \tag{8}$$

**Remark 1.** In contrast to the error bound for final estimators  $\hat{U}_k$  in (7), an intermediate step in the proof for Theorem 3 yields the following upper bound for initializations  $\hat{U}_k^{(0)}$ , i.e., the output from Algorithm 1 Step 1,

$$\mathbb{E}r_k^{-1/q} \left\| \sin \Theta \left( \hat{U}_k^{(0)}, U_k^{(0)} \right) \right\|_q \le C \frac{\sqrt{p_k}}{\lambda/\sigma} + \frac{Cp^{3/2}}{\lambda^2/\sigma^2}, \quad k = 1, 2, 3.$$
 (9)

Compared to Theorem 1, the bound in (9) is suboptimal as long as  $\lambda/\sigma = p^{\alpha}$  when  $3/4 \le \alpha < 1$ . Thus, the higher-order SVD (HOSVD)  $\hat{U}_{k}^{(0)}$  [28] may yield sub-optimal result. We will further illustrate this phenomenon by numerical analysis in Section 5.

Remark 2. Especially when r=1, Theorem 1 confirms the heuristic conjecture raised in Richard and Montanari [23] that the tensor unfolding method yields reliable estimates for order-3 spiked tensors provided that  $\lambda/\sigma > O(p^{3/4})$ . Moreover, Theorem 1 further shows the power iterations are necessary in order to refine the reliable estimates to minimax-optimal estimates.

Our result in Theorem 1 outperforms the ones by Sum-of-Squares (SOS) scheme (see, e.g., [24, 36]), where an additional logarithm factor on the assumption of  $\lambda$  is required. In addition, the method we analyze here, i.e., HOOI, is efficient, easy to implement, and achieves optimal rate of convergence for estimation error.

Remark 3. The strong SNR assumption  $(\lambda/\sigma \ge Cp^{3/4})$  is crucial to guarantee the performance of Algorithm 1. Actually, to ensure that Step 1 in Algorithm 1 provides meaningful initializations,  $\lambda$  should be at least of order  $p^{3/4}$  according to our theoretical analysis.

Moreover, the estimators with high likelihood, such as MLE, achieve the following upper bounds under weaker assumption that  $\lambda/\sigma \geq Cp^{1/2}$ .

**Theorem 2** (Upper Bound for Estimators with Large Likilihood and MLE). Suppose there exist constants  $C_0, c_0 > 0$  such that  $p_k \leq C_0 p$ ,  $r_k \leq C_0 p^{1/2}$ ,  $\|\mathbf{X}\|_F \leq C_0 \sigma\left(\exp(c_0 p)\right)$ ,  $\max\{r_1, r_2, r_3\} \leq C_0 \min\{r_1, r_2, r_3\}$  for  $p = \min\{p_1, p_2, p_3\}$  and k = 1, 2, 3. Suppose  $\hat{U}_k^{\bullet} \in \mathbb{O}_{p_k, r_k}$  are estimators satisfying

$$\min_{\hat{\mathbf{S}}^{\bullet}} \|\hat{\mathbf{Y}} - \hat{\mathbf{S}}^{\bullet} \times_{1} \hat{U}_{1}^{\bullet} \times \hat{U}_{2}^{\bullet} \times_{3} \hat{U}_{3}^{\bullet} \|_{F}^{2} \leq \min_{\hat{\mathbf{S}}} \|\hat{\mathbf{Y}} - \hat{\mathbf{S}} \times_{1} U_{1} \times U_{2} \times_{3} U_{3} \|_{F}^{2}, \tag{10}$$

i.e., the likelihood value of  $\hat{U}_k^{\bullet}$  is no less than  $U_k$ . Then there exists a uniform constant  $C_{gap} > 0$  (which does not depend on  $p_k, r_k, \lambda, \sigma, q$ ) such that whenever

$$\lambda/\sigma \geq C_{gap}p^{1/2}$$
, (i.e., in moderate or strong SNR cases),

$$\hat{U}_{1}^{\bullet}, \hat{U}_{2}^{\bullet}, \hat{U}_{3}^{\bullet}, \text{ and } \hat{\mathbf{X}}^{\bullet} = \hat{\mathbf{S}}^{\bullet} \times_{1} \hat{U}_{1}^{\bullet} \times_{2} \hat{U}_{2}^{\bullet} \times_{3} \hat{U}_{3}^{\bullet} \text{ satisfy}$$

$$\mathbb{E}r_{k}^{1/q} \left\| \sin \Theta(\hat{U}_{k}^{\bullet}, U_{k}) \right\|_{q} \leq \frac{C\sqrt{p_{k}}}{\lambda/\sigma}, \quad k = 1, 2, 3, \quad 1 \leq q \leq 2,$$

$$\mathbb{E} \left\| \hat{\mathbf{X}}^{\bullet} - \mathbf{X} \right\|_{F}^{2} \leq C\sigma^{2} \left( p_{1}r_{1} + p_{2}r_{2} + p_{3}r_{3} \right),$$

$$\mathbb{E} \frac{\|\hat{\mathbf{X}}^{\bullet} - \mathbf{X}\|_{F}^{2}}{\|\mathbf{X}\|_{F}^{2}} \leq C \left( \frac{p_{1} + p_{2} + p_{3}}{\lambda^{2}/\sigma^{2}} \bigwedge 1 \right).$$
(11)

Especially, the upper bounds of (11) hold for maximum likelihood estimators (4).

Then we establish the lower bound for tensor SVD. We especially consider the following class of general low-rank tensors,

$$\mathcal{F}_{p,r}(\lambda) = \left\{ \mathbf{X} \in \mathbb{R}^{p_1 \times p_2 \times p_3} : \operatorname{rank}_k(\mathbf{X}) \le r_k, \sigma_{r_k} \left( \mathcal{M}_k(\mathbf{X}) \right) \ge \lambda, k = 1, 2, 3 \right\}. \tag{12}$$

Here  $\mathbf{p} = (p_1, p_2, p_3)$ ,  $\mathbf{r} = (r_1, r_2, r_3)$  represent the dimension and rank triplets,  $\lambda$  is the smallest non-zero singular value for each matricization of  $\mathbf{X}$ , which essentially measures the signal strength of the problem. The following lower bound holds over  $\mathcal{F}_{\mathbf{p},\mathbf{r}}(\lambda)$ .

**Theorem 3** (Lower Bound). Suppose  $p = \min\{p_1, p_2, p_3\}$ ,  $\max\{p_1, p_2, p_3\} \le C_0 p$ ,  $\max\{r_1, r_2, r_3\} \le C_0 \min\{r_1, r_2, r_3\}$ ,  $4r_1 \le r_2 r_3$ ,  $4r_2 \le r_3 r_1$ ,  $4r_3 \le r_1 r_2$ ,  $1 \le r_k \le p_k/3$ , and  $\lambda > 0$ , then there exists a universal constant c > 0 such that for  $1 \le q \le \infty$ ,

$$\inf_{\tilde{U}_k} \sup_{\mathbf{X} \in \mathcal{F}_{\boldsymbol{p}, \boldsymbol{r}}(\lambda)} \mathbb{E} r_k^{-1/q} \left\| \sin \Theta \left( \tilde{U}_k, U_k \right) \right\|_q \ge c \left( \frac{\sqrt{p_k}}{\lambda / \sigma} \bigwedge 1 \right), \quad k = 1, 2, 3, \tag{13}$$

$$\inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X} \in \mathcal{F}_{\boldsymbol{p}, \boldsymbol{r}}(\lambda)} \mathbb{E} \left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_{\mathrm{F}}^{2} \ge c\sigma^{2} \left( p_{1}r_{1} + p_{2}r_{2} + p_{3}r_{3} \right),$$

$$\inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X} \in \mathcal{F}_{\boldsymbol{p}, \boldsymbol{r}}(\lambda)} \mathbb{E} \frac{\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}}^{2}}{\|\mathbf{X}\|_{\mathrm{F}}^{2}} \ge c \left( \frac{p_{1} + p_{2} + p_{3}}{\lambda^{2}/\sigma^{2}} \bigwedge 1 \right).$$
(14)

**Remark 4.** Theorem 3 contains two folds of meanings. First, when  $\lambda/\sigma \leq cp^{1/2}$  for some small constant c > 0, i.e., under weak SNR setting, the constant term dominates in (13), and there are no consistent estimates for  $U_1, U_2, U_3$ . Secondly, when  $\lambda/\sigma \geq Cp^{1/2}$ , i.e., under strong and moderate SNR settings,  $\frac{\sqrt{p_k}}{\lambda/\sigma}$  dominates in (13), which provides non-trivial statistical lower bounds for estimation errors.

We further define  $\tau^2 = \mathbb{E} \|\mathbf{Z}\|_F^2$  as the expected squared Frobenius norm of the whole tensor. In summary, Theorems 1, 2, and 3 together yield the following statistical limits for tensor SVD.

1. Under strong SNR that  $\lambda/\sigma \geq Cp^{3/4}$  (or  $\lambda/\tau \geq Cp^{-3/4}$ ), the higher-order orthogonal iteration, i.e., Algorithm 1, provides minimax rate-optimal estimators for  $U_1, U_2, U_3$ , and  $\mathbf{X}$ .

$$\inf_{\hat{\mathbf{U}}_{k}} \sup_{\mathbf{X} \in \mathcal{F}_{\boldsymbol{p},\boldsymbol{r}}(\lambda)} \mathbb{E}r_{k}^{-1/q} \left\| \sin \Theta(\hat{U}_{k}, U_{k}) \right\|_{q} \approx \frac{\sqrt{p_{k}}}{\lambda/\sigma} \approx \frac{\tau\sqrt{p_{k}}}{\lambda\sqrt{p_{1}p_{2}p_{3}}}, \quad k = 1, 2, 3, 1 \leq q \leq +\infty, \\
\inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X} \in \mathcal{F}_{\boldsymbol{p},\boldsymbol{r}}(\lambda)} \mathbb{E} \left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_{F}^{2} \approx \sigma^{2} \left( p_{1}r_{1} + p_{2}r_{2} + p_{3}r_{3} \right) \approx \frac{\tau^{2}(p_{1}r_{1} + p_{2}r_{2} + p_{3}r_{3})}{p_{1}p_{2}p_{3}}, \\
\inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X} \in \mathcal{F}_{\boldsymbol{p},\boldsymbol{r}}(\lambda)} \mathbb{E} \frac{\|\hat{\mathbf{X}} - \mathbf{X}\|_{F}^{2}}{\|\mathbf{X}\|_{F}^{2}} \approx \left( \frac{p_{1} + p_{2} + p_{3}}{\lambda^{2}/\sigma^{2}} \wedge 1 \right) \approx \left( \frac{\tau^{2}(p_{1} + p_{2} + p_{3})}{\lambda^{2}p_{1}p_{2}p_{3}} \wedge 1 \right). \tag{15}$$

- 2. Under moderate SNR that  $Cp^{1/2} \leq \lambda/\sigma \leq cp^{3/4}$  (or  $Cp^{-1} \leq \lambda/\tau \leq cp^{-3/4}$ ), the estimators with high likelihood (10), including the MLE (4), are minimax rate-optimal. The rate here is exactly the same as (15).
- 3. Under weak SNR that  $\lambda/\sigma \leq cp^{1/2}$  (or  $\lambda/\tau \leq cp^{-1}$ ), there are no consistent estimators for  $U_1, U_2, U_3$ , or **X**.

However, as we have discussed in Section 2.2, MLE is not applicable even with moderate dimension. It is still crucial to know whether there is any fast and efficient algorithm for tensor SVD under moderate SNR setting.

## 4 Computational Limits in Moderate SNR Case

In this section, we focus on the computational aspect of tensor SVD under moderate SNR setting. If  $\lambda/\sigma = p^{\alpha}$  with  $p = \min\{p_1, p_2, p_3\}$ ,  $\alpha < 3/4$ , we particularly develop the computational lower bound to show that every polynomial-time algorithm is statistically inconsistent in estimating  $U_1$ ,  $U_2$ ,  $U_3$ , and  $\mathbf{X}$  based on computational hardness assumption. In recent literature, we have seen achievements in obtaining computational lower bounds via computational hardness assumptions for many problems, such as sparse PCA [45, 46, 47, 48], submatrix localization [49, 50, 51], tensor completion [52], sparse CCA [48], and community detection [53]. The computational hardness assumptions, such as planted clique detection and Boolean satisfiability, has been widely studied and conjectured that no polynomial-time algorithm exists under certain settings. For tensor SVD, our computational lower bound is established upon the hardness hypothesis of hypergraphic planted clique detection, which is discussed in details in the next section.

#### 4.1 Planted clique detection in hypergraphs

Let G=(V,E) be a graph, where  $V=\{1,2,\ldots,N\}$  and E are the vertex and edge sets, respectively. For a standard graph, the edge  $e=(i,j)\in E$  indicates certain relation exists between vertices i and j in V. A 3-hypergraph (or simply noted as hypergraph, without causing confusion) is a natural extension, where each hyper-edge is represented by an unordered group of three different vertices, say  $e=(i,j,k)\in E$ . Given a hypergraph G=(V,E) with |V|=N, its adjacency tensor  $\mathbf{A}\in\{0,1\}^{N\times N\times N}$  is defined as

$$A_{ijk} = \begin{cases} 1, & \text{if } (i, j, k) \in E; \\ 0, & \text{otherwise.} \end{cases}$$

We denote the Erdős-Rényi hypergraph of N vertices as  $\mathcal{G}_3(N, 1/2)$ , if for each  $1 \leq i < j < k \leq N$ , (i, j, k) is included into the hyper-edge set independently with probability 1/2. For  $V_1 \subset V$  and certain integer  $1 \leq \kappa_N \leq |V_1|$ , we use  $\mathcal{G}_3(N, 1/2, \kappa_N, V_1)$  to denote a random hypergraph where a clique of size  $\kappa_N$  is planted inside  $V_1$ . More precisely, we first sample a random graph from  $\mathcal{G}_3(N, 1/2)$ , then pick  $\kappa_N$  vertices uniformly at random from  $V_1$ , denote them as C, and connecting all hyper-edges (i, j, k) for all distinct triplets  $i, j, k \in C$ . Conventionally, the planted

clique detection is referred to as the problem for distinguishing whether there is any planted clique hidden in the Erdős-Rényi graph. To simplify our analysis in tensor SVD later, we propose a slightly different version of hypergraphic planted clique detection problem as follows.

**Definition 1.** Let G be drawn from either  $\mathcal{G}_3(N, 1/2, \kappa_N, V_1)$  or  $\mathcal{G}_3(N, 1/2, \kappa_N, V_2)$ , where  $V_1 = \{1, 2, ..., \lfloor N/2 \rfloor \}$  and  $V_2 = \{\lfloor N/2 \rfloor + 1, \lfloor N/2 \rfloor + 2, ..., N \}$ . The hypergraphic planted clique detection problem, noted as  $\mathbf{PC}_3(N, \kappa_N)$ , refers to the hypothesis testing problem

$$H_0: G \sim \mathcal{G}_3(N, 1/2, \kappa_N, V_1) \quad vs. \quad H_1: G \sim \mathcal{G}_3(N, 1/2, \kappa_N, V_2).$$
 (16)

Given a hypergraph G sampled from either  $H_0$  or  $H_1$  with adjacency tensor  $\mathbf{A} \in \{0,1\}^{N \times N \times N}$ , let  $\psi(\cdot) : \{0,1\}^{N \times N \times N} \mapsto \{0,1\}$  be a binary-valued function on  $\mathbf{A}$  such that  $\psi(\mathbf{A}) = 1$  indicates rejection of  $H_0$ . Then the risk of test  $\psi$  is defined as the sum of Type-I and II errors,

$$\mathcal{R}_{N,\kappa_N}(\psi) = \mathbb{P}_{H_0}\{\psi(\mathbf{A}) = 1\} + \mathbb{P}_{H_1}\{\psi(\mathbf{A}) = 0\}.$$

Put it differently, given a random hypergraph  $G \sim H_0$  or  $H_1$ , our goal is to identify whether the clique is planted in the first or second half of vertices.

When we replace the hyper-edges (involving three vertices each) of  $\mathcal{G}_3(N, 1/2, \kappa_N, V_1)$  by the regular edges (involving two vertices each), the above hypergraphic planted clique detection becomes the traditional planted clique detection problem. To provide circumstantial evidence to the hardness of  $\mathbf{PC}_3(N, \kappa_N)$ , it is helpful for us to review some well-known results of the traditional planted clique detection here. First, the difficulty of traditional planted clique detection depends crucially on the planted clique size:  $\kappa_N$ . [54] and [55] showed that if  $\kappa_N = o(\log N)$ , it is statistically impossible to determine whether a planted clique exists since a random graph  $G \sim \mathcal{G}_2(N, 1/2)$  contains a clique of size  $2 \log N$  with high probability. When  $\kappa_N \geq C\sqrt{N}$ , it has been shown that the planted clique can be located by performing polynomial-time operations by spectral methods [56, 57]. If the size clique further increases, say  $\kappa_N \geq C\sqrt{N\log N}$ , [58] developed an algorithm to find exactly the planted clique with high probability in polynomial time. However, when  $\log N \ll \kappa_N \ll \sqrt{N}$ , there is still no known polynomial-time algorithm for planted clique detection, and it is currently widely conjectured by the theoretical computer science and graph theory community that such polynomial-time algorithm may not exist (see [59], [60], [55] and references therein).

When moving to hypergraphs, the hardness of  $\mathbf{PC}_3(N, \kappa_N)$ , to the best of our knowledge, remains unclear. In an extreme case of exhaustive search, it needs an exponential number of operations, i.e.,  $\binom{N}{\kappa_N}$ , to verify a solution. In addition, the performance of the simple matricization-spectral method (which shares similar idea as the proposed Algorithm 1) highly depends on the size of the clique  $\kappa_N$ . We particularly have the following Proposition 1.

**Proposition 1.** Suppose  $G \sim \mathcal{G}_3(N, 1/2, \kappa_N, V_1)$ , so there exists  $C \subseteq V_1$  as a planted clique of size  $\kappa_N$  with uniform random position. Let  $\mathbf{A}$  be the corresponding adjacency tensor, and  $1_C \in \mathbb{R}^{|V_1|}$  be the indicator for the hidden clique that  $(1_C)_i = 1_{\{i \in C\}}$ . We further partition  $V_1 = \{1, \ldots, \lfloor N/2 \rfloor\}$  into three equal subsets:  $D_k = \{\lfloor kN/6 \rfloor + 1, \ldots, \lfloor (k+1)N/6 \rfloor\}$  for k = 1, 2, 3. Then we can calculate  $\hat{u}_k \in \mathbb{R}^{|V_1|}$  as the leading left singular vector of  $\mathcal{M}_k(2 \cdot \mathbf{A}_{[D_1, D_2, D_3]} - 1_{|D_1| \times |D_2| \times |D_3|})$ , where  $1_{|D_1| \times |D_2| \times |D_3|}$  is a  $|D_1|$ -by- $|D_2|$ -by- $|D_3|$  tensor with all entries 1. If the sequence  $\{\kappa_N\}$  satisfies  $\liminf_{N \to \infty} \frac{\kappa_N}{N^{1/2}} = \infty$ , then

$$\sin\Theta(\hat{u}_k,(1_C)_{D_k}) \stackrel{d}{\to} 0$$
, as  $N \to \infty$ ,  $k = 1,2,3$ .

In another word, the angle between  $\hat{u}_k$  and  $(1_C)_{D_k}$  tends to 0 in probability.

**Remark 5.** For technical convenience, we partition  $V_1$  into three parts and perform SVD on  $\mathcal{M}_k(2\mathbf{A}_{[D_1,D_2,D_3]}-1_{|D_1|\times|D_2|\times|D_3|})$  to ensure that most of the entries of  $\mathcal{M}_k(\mathbf{A})$  are i.i.d. Rademacher distributed.

Proposition 1 suggests that  $\hat{u}_k$  can be used to locate C when  $\kappa_N \gg N^{1/2}$ . However, the theoretical analysis in Proposition 1 fails when  $\kappa_N = N^{(1-\tau)/2}$  for  $\tau > 0$ , and we conjecture that such computational barrier is essential. Particularly, we propose the following computational hardness assumption on hypergraphic planted clique detection.

**Hypothesis 1. H(\tau).** For any sequence  $\{\kappa_N\}$  such that  $\lim_{N\to\infty} \sup \frac{\log \kappa_N}{\log \sqrt{N}} \leq (1-\tau)$  and any sequence of polynomial-time tests  $\{\psi_N\}$ ,

$$\liminf_{N\to\infty} \, \mathcal{R}_{N,\kappa_N}(\psi_N) \ge \frac{1}{2}.$$

#### 4.2 The computational lower bound of tensor SVD

Now we are ready to develop the computational lower bound for tensor SVD based on Hypothesis  $\mathbf{H}(\tau)$ . Recall

$$\mathbf{Y} = \mathbf{X} + \mathbf{Z} \in \mathbb{R}^{p_1 \times p_2 \times p_3}, \quad \mathbf{X} = \mathbf{S} \times_1 U_1 \times_2 U_2 \times_3 U_3, \quad \mathbf{Z} \stackrel{iid}{\sim} N(0, \sigma^2).$$

To better present the asymptotic argument, we add a superscript of dimension,  $p = \min\{p_1, p_2, p_3\}$ , to the estimators, i.e.,  $\hat{U}_k^{(p)}$ ,  $\hat{\mathbf{X}}^{(p)}$ . The computational lower bound is then presented as below.

**Theorem 4** (Computational Lower Bound). Suppose hypergraphic planted clique assumption  $\mathbf{H}(\tau)$  holds for some  $\tau \in (0,1)$ . Then there exist absolute constants  $c_0, c_1 > 0$  such that if  $\lambda/\sigma \leq c_0 \left(\frac{p^{3(1-\tau)/4}}{\sqrt{\log 3p}}\right)$ , for any integers  $r_1, r_2, r_3 \geq 1$  and any polynomial time estimators  $\hat{U}_k^{(p)}$ ,  $\hat{\mathbf{X}}^{(p)}$ , the following inequalities hold

$$\lim_{p \to \infty} \inf_{\mathbf{X} \in \mathcal{F}_{p,r}(\lambda)} \mathbb{E} \left\| \sin \Theta \left( \hat{U}_k^{(p)}, U_k \right) \right\|^2 \ge c_1, \quad k = 1, 2, 3,$$
(17)

$$\liminf_{p \to \infty} \sup_{\mathbf{X} \in \mathcal{F}_{\mathbf{p}, \mathbf{r}}(\lambda)} \frac{\mathbb{E} \|\hat{\mathbf{X}}^{(p)} - \mathbf{X}\|_{\mathrm{F}}^2}{\|\mathbf{X}\|_{\mathrm{F}}^2} \ge c_1.$$
(18)

**Remark 6.** For technical reasons, there is an additional logarithmic factor in the condition  $\lambda/\sigma \leq c_0 \left(\frac{p^{3(1-\tau)/4}}{\sqrt{\log 3p}}\right)$ , compared with the statistical lower bound in Theorem 3. Since  $\tau$  is a strictly positive number, the effect of logarithmic factor is dominated by  $p^c$  for any c > 0 asymptotically.

Theorem 4 illustrates the computational hardness for tensor SVD under moderate scenario that  $\lambda/\sigma = p^{\alpha}, 1/2 \le \alpha < 3/4$ , if the hypergraphic planted clique assumption  $\mathbf{H}(\tau)$  holds for any  $\tau > 0$ .

## 5 Simulations

In this section, we further illustrate the statistical and computational limits for tensor SVD via numerical studies.

We first consider the average Shatten q-sin  $\Theta$ -norm losses for initial estimators  $\hat{U}_k^{(0)}$  (HOSVD) and final estimators  $\hat{U}_k$  (HOOI) under the following simulation setting. For given  $(p, r, \lambda)$ , we

| $(p,r,\lambda)$ | $l_1(\hat{U})$ | $l_1(\hat{U}^{(0)})$ | $l_2(\hat{U})$ | $l_2(\hat{U}^{(0)})$ | $l_5(\hat{U})$ | $l_5(\hat{U}^{(0)})$ | $l_{\infty}(\hat{U})$ | $l_{\infty}(\hat{U}^{(0)})$ |
|-----------------|----------------|----------------------|----------------|----------------------|----------------|----------------------|-----------------------|-----------------------------|
| (50, 5, 20)     | 1.1094         | 2.1192               | 0.5194         | 1.0535               | 0.3572         | 0.7991               | 0.3286                | 0.7699                      |
| (50, 5, 50)     | 0.4297         | 0.5243               | 0.2016         | 0.2519               | 0.1392         | 0.1815               | 0.1283                | 0.1713                      |
| (50, 10, 20)    | 2.4529         | 4.5208               | 0.8179         | 1.5674               | 0.4629         | 0.9611               | 0.3955                | 0.8762                      |
| (50, 10, 50)    | 0.9111         | 1.1210               | 0.3030         | 0.3771               | 0.1707         | 0.2175               | 0.1452                | 0.1890                      |
| (100, 5, 40)    | 0.7952         | 1.5649               | 0.3695         | 0.7707               | 0.2509         | 0.5778               | 0.2294                | 0.5543                      |
| (100, 5, 60)    | 0.5301         | 0.8132               | 0.2463         | 0.3938               | 0.1673         | 0.2878               | 0.1530                | 0.2731                      |
| (100, 10, 40)   | 1.7448         | 3.5371               | 0.5688         | 1.1943               | 0.3087         | 0.7015               | 0.2554                | 0.6246                      |
| (100, 10, 60)   | 1.1466         | 1.8055               | 0.3735         | 0.6015               | 0.2021         | 0.3427               | 0.1660                | 0.2975                      |

Table 1: The average Schatten q-sin  $\Theta$  loss of the final estimations  $\hat{U}_k$  and the spectral initializations  $\hat{U}_k^{(0)}$  based on 100 repetitions. Here,  $p_1 = p_2 = p_3 = p$ ,  $r_1 = r_2 = r_3 = r$ ,  $l_q(\hat{U}) = \frac{1}{3} \sum_{k=1}^3 \|\sin\Theta(\hat{U}_k, U_k)\|_q$ .

let  $p=p_1=p_2=p_3, r=r_1=r_2=r_3$ , generate  $\tilde{U}_k\in\mathbb{R}^{p_k\times r_k}$  as i.i.d. standard Gaussian matrices, and apply QR decomposition on  $\tilde{U}_k$  and assign the Q part to  $U_k$ . In other words, the singular subspaces  $U_1, U_2, U_3$  are drawn randomly from Haar measure. Then we construct  $\tilde{\mathbf{S}}\in\mathbb{R}^{r_1\times r_2\times r_3}$  as an i.i.d. Gaussian tensor, and rescale as  $\mathbf{S}=\tilde{\mathbf{S}}\cdot\frac{\lambda}{\min_{k=1,2,3}\sigma_{r_k}(\mathcal{M}_k(\tilde{\mathbf{S}}))}$  to ensure  $\min_{k=1,2,3}\sigma_{r_k}(\mathcal{M}_k(\mathbf{X}))\geq\lambda$ . Next, we construct  $\mathbf{Y}=\mathbf{X}+\mathbf{Z}$ , where the signal tensor  $\mathbf{X}=\mathbf{S}\times_1 U_1\times_2 U_2\times_3 U_3$ , the noise tensor  $\mathbf{Z}$  are drawn from i.i.d. standard Gaussian distribution. We apply Algorithm 1 to  $\mathbf{Y}$ , and record the average numerical performance under various values of  $p,r,\lambda$ . The results based on 100 replications are shown in Table 1. We can clearly see that the power iterations (Step 2 in Algorithm 1, i.e., HOOI) significantly improve upon spectral initializations (Step 1 in Algorithm 1, i.e., HOSVD) in different Shatten-q sin  $\Theta$  losses under various settings.

Then we consider another setting that  $\mathbf{X}$  has three different dimensions. Specifically, we generate  $\mathbf{Y} = \mathbf{X} + \mathbf{Z}$  by the same scheme as the previous setting with varying  $(p_1, p_2, p_3)$  and fixed  $r_1 = r_2 = r_3 = 5$ . We repeat the experiment for 100 times, then record the average estimation errors in Table 2. Again, we can see HOOI performs well under various values of dimensions.

| $(p_1, p_2, p_3, \lambda)$ | $l_{\infty}(\hat{U}_1)$ | $l_2(\hat{U}_1)$ | $l_{\infty}(\hat{U}_2)$ | $l_2(\hat{U}_2)$ | $l_{\infty}(\hat{U}_3)$ | $l_2(\hat{U}_3)$ | $\ \hat{\mathbf{X}} - \mathbf{X}\ _{\mathrm{F}}$ | $\frac{\ \hat{\mathbf{X}} - \mathbf{X}\ _{\mathrm{F}}}{\ \mathbf{X}\ _{\mathrm{F}}}$ |
|----------------------------|-------------------------|------------------|-------------------------|------------------|-------------------------|------------------|--|--|
| (20, 30, 50, 20)           | 0.2082                  | 0.3032           | 0.2530                  | 0.3858           | 0.3109                  | 0.4975           | 24.7037  | 0.3276   |
| (20, 30, 50, 100)          | 0.0409                  | 0.0596           | 0.0498                  | 0.0761           | 0.0641                  | 0.1017           | 23.5708  | 0.0631   |
| (30, 50, 100, 20)          | 0.2674                  | 0.4036           | 0.3354                  | 0.5247           | 0.4456                  | 0.7252           | 33.6219  | 0.4479   |
| (30, 50, 100, 100)         | 0.0490                  | 0.0753           | 0.0640                  | 0.1012           | 0.0911                  | 0.1469           | 30.9540  | 0.0822   |
| (100, 200, 300, 50)        | 0.1840                  | 0.2982           | 0.2551                  | 0.4301           | 0.3161                  | 0.5155           | 57.8482  | 0.3090   |
| (100, 200, 300, 100)       | 0.0940                  | 0.1506           | 0.1259                  | 0.2117           | 0.1638                  | 0.2627           | 55.9009  | 0.1505   |
| (200, 300, 400, 50)        | 0.2579                  | 0.4335           | 0.3331                  | 0.5523           | 0.3420                  | 0.6017           | 72.2912  | 0.4026   |
| (200, 300, 400, 150)       | 0.0825                  | 0.1389           | 0.1076                  | 0.1739           | 0.1277                  | 0.2024           | 68.0305  | 0.1199   |

Table 2: The average spectral and Frobenius  $\sin\Theta$  loss for  $\hat{U}_1$ ,  $\hat{U}_2$ ,  $\hat{U}_3$  and average Frobenius loss for  $\hat{\mathbf{X}}$  under various settings. Here  $l_{\infty}(\hat{U}_k) = \|\sin\Theta(\hat{U}_k, U_k)\|$ ,  $l_2(\hat{U}) = \|\sin\Theta(\hat{U}_k, U_k)\|_{\mathrm{F}}$ .

Next, we illustrate the phase transition phenomenon of tensor SVD. Let  $\mathbf{X} = \mathbf{S} \times_1 U_1 \times_2 U_2 \times_3 U_3$  be a p-by-p-by-p tensor, where  $U_1, U_2, U_3$  are randomly generated p-by-r orthogonal matrices from Haar measure and  $\mathbf{S} \in \mathbb{R}^{r \times r \times r}$  is a fixed diagonal tensor such that  $S_{i,j,k} = p^{\alpha} \cdot 1_{\{i=j=k\}}, 1 \le i, j, k \le r$  for  $\alpha \in [0.4, 0.9]$ . Then  $p^{\alpha}$  is the signal strength in our context. The entries of  $\mathbf{Z}$  are generated as either i.i.d. N(0,1) or  $\text{Unif}[-\sqrt{3},\sqrt{3}]$ , which are sub-Gaussian, mean 0, and variance 1. To demonstrate the phase transitions at both  $p^{3/4}$  and  $p^{1/2}$ , ideally one wishes to implement both MLE and HOOI. Since MLE, i.e., the best low-rank approximation estimator (4), is computationally intractable, we instead consider the following oracle warm-start HOOI to obtain an approximation for MLE: suppose an oracle provides a warm start as

$$\hat{U}_k^{(0)\text{warm}} = \frac{1}{\sqrt{2}}U_k + \frac{1}{\sqrt{2}}U_k', \quad k = 1, 2, 3,$$

where  $U_k$  is the true underlying loading and  $U_k'$  is a p-by-r random orthonormal matrix in the complementary space of  $U_k$ .  $\{U_k'\}_{k=1}^3$  here are generated based on the following scheme: first calculate  $U_{k\perp} \in \mathbb{O}_{p,p-r}$  as the orthogonal complement of  $U_k$ , then construct  $U_k' = U_{k\perp}O$  for some random orthogonal matrix  $O \in \mathbb{O}_{p-r,r}$ . Based on the oracle warm-start, we apply Steps 2 and 3 of Algorithm 1 to obtain the warm-start HOOI estimator  $\hat{U}_k^{\text{warm}}$  as an approximation for MLE.

We let p vary from 50 to 100, r=5, and apply both the spectral-start HOOI (i.e., the original HOOI and Algorithm 1) and the oracle warm-start HOOI. The average spectral  $\sin \Theta$ 

loss, i.e.,  $l_{\infty}(\hat{U}) = \frac{1}{3} \sum_{k=1}^{3} \|\sin\Theta(\hat{U}_{k}, U_{k})\|$  from 100 repetitions are presented in Figure 1, where the upper panel and lower panel correspond to the i.i.d. Gaussian noise and i.i.d. uniform noise cases, respectively. Both panels of Figure 1 clearly demonstrate the phase transition effects: the estimation error significantly decreases around SNR =  $p^{3/4}$  and around SNR =  $p^{1/2}$  for spectral-start HOOI and oracle warm-start HOOI, respectively. This exactly matches our theoretical findings in Section 3. In addition, there is little difference between two plots in the upper and lower panels, which implies that the statistical estimation error for tensor SVD mainly relies on the SNR and is less influenced by the particular sub-Gaussian noise type.

#### 6 Discussions: Further Generalizations

In this article, we propose a general framework for tensor singular value decomposition (tensor SVD), which focuses on extracting the underlying Tucker low-rank structure from the noisy tensor observations. We provide a comprehensive analysis on tensor SVD in aspects of both statistics and computation. The problem exhibits three distinct phases according to the signal-to-noise ratio (SNR): with strong SNR, the higher order orthogonal iteration (HOOI) performs efficiently and achieves statistical optimal results; with weak SNR, no method performs consistently; with moderate SNR, the estimators with high likelihood, such as the computational intractable MLE, perform optimally in statistical convergence rate, and no polynomial algorithm can do so unless we have a polynomial-time algorithm for the hypergraphic planted clique problem.

The results of this paper are mainly presented under i.i.d. Gaussian noise. However, when the noise is more generally i.i.d. sub-Gaussian distributed, say

$$\mathbf{Z} \stackrel{iid}{\sim} Z$$
, where  $\|Z\|_{\psi_2} = \sup_{q \ge 1} q^{-1/2} (\mathbb{E}|Z|^q)^{1/q} \le \sigma$ ,

we can derive the upper bound results similarly as Theorems 1 and 2, as the proofs of main technical tools, including Lemmas 5 and 8, still hold for i.i.d. sub-Gaussian noise case.

We have also focused our presentations mainly on order-3 tensors throughout this article. The results can be additionally generalized to order-d tensor SVD for any  $d \ge 2$ . Suppose one

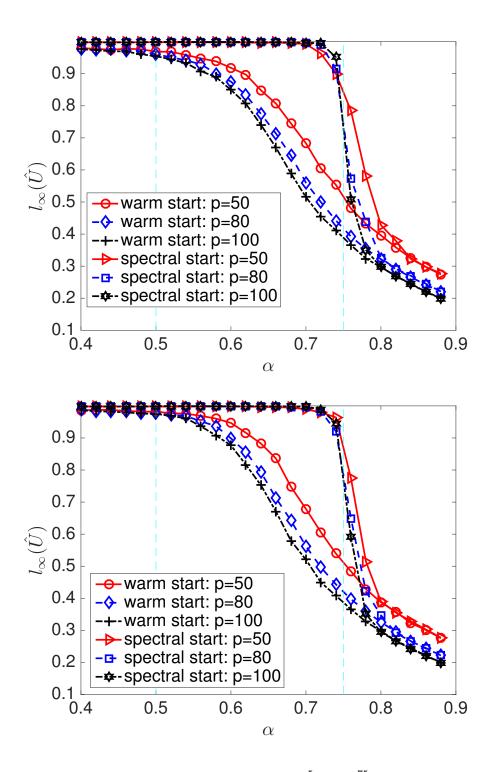


Figure 1: Phase transitions in tensor SVD at SNR =  $p^{.5}$  and  $p^{.75}$ . Upper panel: Gaussian noise N(0,1). Lower panel: uniform noise Unif $[-\sqrt{3},\sqrt{3}]$ .

observes an order-d tensor as follows,

$$\mathbf{Y} = \mathbf{X} + \mathbf{Z}, \quad \mathbf{X} = \mathbf{S} \times_1 U_1 \cdots \times_d U_d,$$
 (19)

where  $\mathbf{Y}, \mathbf{X}, \mathbf{Z} \in \mathbb{R}^{p_1 \times \cdots \times p_d}$ ,  $U_k \in \mathbb{O}_{p_k, r_k}$  for  $k = 1, \dots, d$ , and  $\mathbf{S} \in \mathbb{R}^{r_1 \times \cdots \times r_d}$ . The higher-order orthogonal iteration can be similarly written as follows [29],

Step 1 Initialize by singular value decomposition of each matricizations,

$$\hat{U}_k^{(0)} = \text{SVD}_{r_k} \left( \mathcal{M}_k(\mathbf{Y}) \right), \quad k = 1, \dots, d.$$

Step 2 Let t = 0, 1, ..., update the estimates for  $U_k$  sequentially for k = 1, ..., d,

$$\hat{U}_{k}^{(t+1)} = \text{SVD}_{r_{k}} \left( \mathbf{Y} \times_{1} (\hat{U}_{1}^{(t+1)})^{\top} \times \dots \times_{(k-1)} (\hat{U}_{k-1}^{(t+1)})^{\top} \times \dots \times_{(k-1)} (\hat{U}_{k-1}^{(t+1)})^{\top} \times \dots \times_{d} (\hat{U}_{d}^{(t)})^{\top} \right).$$

The iteration is continued until convergence or maximum number of iteration is reached. Step 3 With the final estimators  $\{\hat{U}_k\}_{k=1}^d$  from Step 2, one estimates **X** as

$$\hat{\mathbf{X}} = \mathbf{Y} \times_1 P_{\hat{U}_1} \times \cdots \times_d P_{\hat{U}_d}.$$

Meanwhile, the non-convex maximum likelihood estimates can be written as

$$(\hat{U}_{1}^{\text{mle}}, \dots, \hat{U}_{d}^{\text{mle}}) = \underset{V_{k} \in \mathbb{O}_{p_{k}, r_{k}}}{\operatorname{arg max}} \left\| \mathbf{Y} \times_{1} V_{1}^{\top} \times \dots \times_{d} V_{d}^{\top} \right\|_{F},$$

$$k=1,\dots,d$$

$$\hat{\mathbf{X}}^{\text{mle}} = \mathbf{Y} \times_{1} P_{\hat{U}_{1}^{\text{mle}}} \times \dots \times_{d} P_{\hat{U}_{d}^{\text{mle}}}.$$

$$(20)$$

Again, let  $\lambda = \min_{1 \leq k \leq d} \sigma_{r_k}(\mathcal{M}_k(\mathbf{X}))$  measure the signal strength. For fixed d, when  $p = \min\{p_1, \ldots, p_d\}$ ,  $\max\{p_1, \ldots, p_d\} \leq Cp$ ,  $r_k \leq Cp^{1/(d-1)}$ , similarly as the proofs for Theorems 1, 2, and 3, it is possible to show under strong SNR case, where  $\lambda/\sigma = p^{\alpha}$  for  $\alpha \geq d/4$ , HOOI achieves optimal rate of convergence over the following class of low-rank tensors

$$\mathcal{F}_{p,r}(\lambda) = \left\{ \mathbf{X} \in \mathbb{R}^{p_1 \times \dots \times p_d} : \sigma_{r_k}(\mathcal{M}_k(\mathbf{X})) \ge \lambda, k = 1, \dots, d \right\};$$

under the weak SNR where  $\lambda/\sigma = p^{\alpha}$  for  $\alpha < 1/2$ , it is impossible to generally have consistent estimators for  $U_1, \ldots, U_d$ , or **X**; under the moderate SNR, where  $\lambda/\sigma = p^{\alpha}$  for  $\frac{1}{2} \leq \alpha < \frac{d}{4}$ ,

the estimators with high likelihood, such as MLE, achieves optimal statistical performance; while one can develop computational lower bound with computational hardness assumption of a higher-order hypergraphic planted clique detection problem similarly as Theorem 4. We can also see that the gap between statistical and computational limits vanishes if d = 2. This coincides with the previous results in matrix denoising literature (see, e.g. [7, 8, 26]), where standard singular value decomposition achieves both statistical optimality and computational efficiency.

Additionally, if d grows rather than stays as a fixed constant, the asymptotics of tensor SVD in both statistical and computational aspects will be an interesting future project.

### 7 Proofs

We collect the proofs in this section for the main results in this paper. To be specific, the proof for Theorems 1, 3, and 4 will be presented in Sections 7.1, 7.3, and 7.2, respectively. Proofs for Theorem 2, Proposition 1, and additional technical lemmas are postponed to the supplementary materials.

#### 7.1 Proof of Theorem 1

We first consider the proof for Theorem 1. Throughout the proof, we assume the noise level  $\sigma^2 = 1$  without loss of generality. For convenience, we denote

$$X_1 = \mathcal{M}_1(\mathbf{X}), \quad X_2 = \mathcal{M}_2(\mathbf{X}), \quad X_3 = \mathcal{M}_3(\mathbf{X})$$

as the matricizations of **X**. We also denote  $Y_1, Y_2, Y_3, Z_1, Z_2, Z_3$  in the similar fashion. We also let  $r = \max\{r_1, r_2, r_3\}$ .

We divide the proof into steps.

1. In this first step, we consider the performance of initialization step, we particularly prove that for any small constant  $c_0 > 0$ , there exists large constant  $C_{gap} > 0$  such that whenever  $\lambda \geq C_{gap}p^{3/4}$ , we have

$$\left\| \sin \Theta(\hat{U}_k^{(0)}, U_k) \right\| \le c_0 \left( \frac{\sqrt{p_k} \lambda + (p_1 p_2 p_3)^{1/4}}{\lambda^2} \right)$$
 (21)

with probability at least  $1 - C \exp(-cp)$ . The proof of this step is closely related to the proof for Theorem 3 in [26]. Note that

$$\hat{U}_{1}^{(0)} = \text{SVD}_{r_{1}}(Y_{1}), \quad Y_{1} = X_{1} + Z_{1},$$

where  $X_1$  is a fixed matrix satisfying rank $(X_1) = r_1$ ;  $Z_1 \in \mathbb{R}^{p_1 \times (p_2 p_3)}$ ,

 $\{(Z_1)_{ij}\}_{i,j=1}^{p_1,p_2p_3} \stackrel{iid}{\sim} N(0,1)$ . This shares the same setting as the one in Theorem 3 in [26], if one sets  $p_1, p_2$  in the statement of Theorem 3 in [26] respectively as  $p_2p_3, p_1$  in our context. Thus we can essentially follow their proof. Let  $U_{1\perp}$  be the orthogonal complement of  $U_1$ . Then the Appendix Equations (1.15), (1.16) in [26] yields

$$\mathbb{P}\left(\left\|\sin\Theta(\hat{U}_{1}^{(0)}, U_{1})\right\|^{2} \leq \frac{C(\lambda^{2} + p_{2}p_{3})\|U_{1\perp}Y_{1}P_{U_{1}^{\top}Y_{1}}\|^{2}}{\lambda^{4}}\right) \\
\geq 1 - C\exp\left\{-c\frac{\lambda^{4}}{\lambda^{2} + p_{2}p_{3}}\right\}$$
(22)

and

$$\mathbb{P}\left(\|U_{1\perp}Y_{1}P_{U_{1}^{\top}Y_{1}}\| \geq x\right) 
\leq C \exp\left\{Cp_{1} - c \min\left(x^{2}, x\sqrt{\lambda^{2} + p_{2}p_{3}}\right)\right\} + C \exp\left\{-c(\lambda^{2} + p_{2}p_{3})\right\}$$
(23)

for some uniform constant C, c > 0. Since

$$\lambda \ge C_{gap}p^{3/4} \ge C_{gap}c\left(p_1^{1/2} + (p_1p_2p_3)^{1/4}\right),$$

if we set  $x = C\sqrt{p_1}$ , (23) further leads to

$$\mathbb{P}\left(\|U_{1\perp}Y_1P_{U_1^{\top}Y_1}\| \ge C\sqrt{p_1}\right) \le Ce^{-cp} + C\exp(-c(\lambda^2 + p^2)) \le C\exp(-cp). \tag{24}$$

Combining (24) and (22), we have proved (21) for k = 1. The proof for (21) for k = 2, 3 can be similarly written down.

2. After spectral initialization, we assume the algorithm evolves from t=0 to  $t=t_{\rm max}$ , where  $t_{\rm max} \geq C\left(\log(\frac{p}{\lambda}) \vee 1\right)$ . In this step, we derive the perturbation bounds for  $\hat{U}_1^{(t_{\rm max})}, \hat{U}_2^{(t_{\rm max})}, \hat{U}_3^{(t_{\rm max})}$  under the assumptions that  $\lambda \geq C_{gap}^{3/4}$  for large constant  $C_{gap} > 0$  and the following inequalities all holds,

$$\max \left\{ \left\| \sin \Theta(\hat{U}_{1}^{(0)}, U_{1}) \right\|, \left\| \sin \Theta(\hat{U}_{2}^{(0)}, U_{2}) \right\|, \left\| \sin \Theta(\hat{U}_{3}^{(0)}, U_{3}) \right\| \right\} \le \frac{1}{2}, \tag{25}$$

$$\max_{\substack{V_{2} \in \mathbb{R}^{p_{2} \times r_{2}} \\ V_{3} \in \mathbb{R}^{p_{3} \times r_{3}}}} \frac{\|Z_{1} \cdot (V_{2} \otimes V_{3})\|}{\|V_{2}\| \|V_{3}\|} \leq C_{1} \sqrt{pr}, \quad \max_{\substack{V_{3} \in \mathbb{R}^{p_{3} \times r_{3}} \\ V_{1} \in \mathbb{R}^{p_{1} \times r_{1}}}} \frac{\|Z_{2} \cdot (V_{3} \otimes V_{1})\|}{\|V_{3}\| \|V_{1}\|} \leq C_{1} \sqrt{pr},$$

$$\max_{\substack{V_{1} \in \mathbb{R}^{p_{1} \times r_{1}} \\ V_{2} \in \mathbb{R}^{p_{2} \times r_{2}}}} \frac{\|Z_{3} \cdot (V_{1} \otimes V_{2})\|}{\|V_{1}\| \|V_{2}\|} \leq C_{1} \sqrt{pr},$$

$$(26)$$

$$||Z_1(U_2 \otimes U_3)|| \le C_1 \sqrt{p_1}, ||Z_2(U_3 \otimes U_1)|| \le C_2 \sqrt{p_2}, ||Z_3(U_1 \otimes U_2)|| \le C_2 \sqrt{p_3}.$$
 (27)

Recall here that  $U_1, U_2, U_3$  are left singular subspaces for  $X_1, X_2, X_3$ , respectively. We let  $L_t$  be the spectral  $\sin \Theta$  norm error for  $\hat{U}_k^{(t)}$ ,

$$L_t = \max_{k=1,2,3} \left\| \sin \Theta(\hat{U}_k^{(t)}, U_k) \right\|, \quad t = 0, 1, 2, \dots$$
 (28)

Given (25),  $L_0 \leq \frac{1}{2}$ . Next we aim to prove that for  $t = 0, 1, \ldots$ ,

$$L_{t+1} = \max_{k=1,2,3} \left\| \sin \Theta(\hat{U}_k^{(t+1)}, U_k) \right\| \le \frac{C_1 \sqrt{pr}}{\lambda} L_t + \frac{C_2 \sqrt{p}}{\lambda} \le \frac{1}{2}.$$
 (29)

To show (29), we first focus on the upper bound of  $\|\sin\Theta(\hat{U}_1^{(t+1)}, U_1)\|$  when t = 0. Define the following key components in our analysis as follows,

$$Y_1^{(t)} = \mathcal{M}_1 \left( \mathbf{Y} \times_2 \left( \hat{U}_2^{(t)} \right)^\top \times_3 \left( \hat{U}_3^{(t)} \right)^\top \right) \overset{\text{Lemma } 4}{=} Y_1 \cdot \left( \hat{U}_2^{(t)} \otimes \hat{U}_3^{(t)} \right) \in \mathbb{R}^{p_1 \times r_2 r_3},$$

$$X_1^{(t)} = \mathcal{M}_1 \left( \mathbf{X} \times_2 \left( \hat{U}_2^{(t)} \right)^\top \times_3 \left( \hat{U}_3^{(t)} \right)^\top \right) \overset{\text{Lemma } 4}{=} X_1 \cdot \left( \hat{U}_3^{(t)} \otimes \hat{U}_1^{(t)} \right) \in \mathbb{R}^{p_1 \times r_2 r_3},$$

$$Z_1^{(t)} = \mathcal{M}_1 \left( \mathbf{Z} \times_2 \left( \hat{U}_2^{(t)} \right)^\top \times_3 \left( \hat{U}_3^{(t)} \right)^\top \right) \overset{\text{Lemma } 4}{=} Z_1 \cdot \left( \hat{U}_1^{(t)} \otimes \hat{U}_2^{(t)} \right) \in \mathbb{R}^{p_1 \times r_2 r_3}.$$

By definition, the left and right singular subspaces of  $X_1$  are  $U_1 \in \mathbb{O}_{p_1,r_1}$  and  $U_2 \otimes U_3 \in \mathbb{O}_{p_2p_3,r_2r_3}$ . Then,

$$\sigma_{r_{1}}\left(X_{1}^{(t)}\right) = \sigma_{r_{1}}\left(X_{1} \cdot \left(\hat{U}_{2}^{(t)} \otimes \hat{U}_{3}^{(t)}\right)\right) = \sigma_{r_{1}}\left(X_{1} \cdot P_{U_{2} \otimes U_{3}} \cdot \left(\hat{U}_{2}^{(t)} \otimes \hat{U}_{3}^{(t)}\right)\right)$$

$$= \sigma_{r_{1}}\left(X_{1} \cdot (U_{2} \otimes U_{3}) \cdot (U_{2} \otimes U_{3})^{\top} \cdot \left(\hat{U}_{2}^{(t)} \otimes \hat{U}_{3}^{(t)}\right)\right)$$

$$\geq \sigma_{r_{1}}\left(X_{1} \cdot (U_{2} \otimes U_{3})\right) \cdot \sigma_{\min}\left(\left(U_{2} \otimes U_{3}\right)^{\top} \cdot \left(\hat{U}_{2}^{(t)} \otimes \hat{U}_{3}^{(t)}\right)\right)$$

$$= \sigma_{r_{1}}(X_{1}) \cdot \sigma_{\min}\left(\left(U_{2}^{\top}\hat{U}_{2}^{(t)}\right) \otimes \left(U_{3}^{\top}\hat{U}_{3}^{(t)}\right)\right)$$

$$\geq \sigma_{r_{1}}(X_{1}) \cdot \sigma_{\min}\left(U_{2}^{\top}\hat{U}_{2}^{(t)}\right) \cdot \sigma_{\min}\left(U_{3}^{\top}\hat{U}_{3}^{(t)}\right)$$

$$\geq \lambda \cdot \left(1 - L_{t}^{2}\right) \quad \text{(by (28) and Lemma 1 in [26])}.$$

Meanwhile,

$$\begin{aligned} & \left\| Z_{1}^{(t)} \right\| = \left\| Z_{1} \left( \hat{U}_{2}^{(t)} \otimes \hat{U}_{3}^{(t)} \right) \right\| \\ &= \left\| Z_{1} \left( P_{U_{2} \otimes U_{3}} + P_{U_{2} \perp} \otimes U_{3} + P_{I_{p} \otimes U_{3} \perp} \right) \left( \hat{U}_{2}^{(t)} \otimes \hat{U}_{3}^{(t)} \right) \right\| \quad \text{(by Lemma 4)} \\ &\leq \left\| Z_{1} \left( P_{U_{2} \otimes U_{3}} \right) \left( \hat{U}_{2}^{(t)} \otimes \hat{U}_{3}^{(t)} \right) \right\| + \left\| Z_{1} \left( P_{U_{2} \perp} \otimes U_{3} \right) \left( \hat{U}_{2}^{(t)} \otimes \hat{U}_{3}^{(t)} \right) \right\| \\ &+ \left\| Z_{1} \left( P_{U_{2} \otimes U_{3} \perp} \right) \left( \hat{U}_{2}^{(t)} \otimes \hat{U}_{3}^{(t)} \right) \right\| + \left\| Z_{1} \left( P_{U_{2} \perp} \otimes U_{3} \right) \left( \hat{U}_{2}^{(t)} \otimes \hat{U}_{3}^{(t)} \right) \right\| \\ &= \left\| Z_{1} (U_{2} \otimes U_{3}) \left( U_{2} \otimes U_{3} \right)^{\top} \left( \hat{U}_{2}^{(t)} \otimes \hat{U}_{3}^{(t)} \right) \right\| \\ &+ \left\| Z_{1} \left( \left( P_{U_{2} \perp} \hat{U}_{2}^{(t)} \right) \otimes \left( P_{U_{3} \perp} \hat{U}_{3}^{(t)} \right) \right) \right\| \\ &+ \left\| Z_{1} \left( \left( P_{U_{2} \perp} \hat{U}_{2}^{(t)} \right) \otimes \left( P_{U_{3} \perp} \hat{U}_{3}^{(t)} \right) \right) \right\| \\ &+ \left\| Z_{1} \left( \left( P_{U_{2} \perp} \hat{U}_{2}^{(t)} \right) \otimes \left( P_{U_{3} \perp} \hat{U}_{3}^{(t)} \right) \right) \right\| \\ &+ \left\| Z_{1} \left( \left( P_{U_{2} \perp} \hat{U}_{2}^{(t)} \right) \otimes \left( P_{U_{3} \perp} \hat{U}_{3}^{(t)} \right) \right) \right\| \\ &+ \left\| Z_{1} \left( \left( P_{U_{2} \perp} \hat{U}_{2}^{(t)} \right) \otimes \left( P_{U_{3} \perp} \hat{U}_{3}^{(t)} \right) \right) \right\| \\ &+ \left\| Z_{1} \left( \left( P_{U_{2} \perp} \hat{U}_{2}^{(t)} \right) \otimes \left( P_{U_{3} \perp} \hat{U}_{3}^{(t)} \right) \right\| \\ &+ \left\| Z_{1} \left( \left( P_{U_{2} \perp} \hat{U}_{2}^{(t)} \right) \otimes \left( P_{U_{3} \perp} \hat{U}_{3}^{(t)} \right) \right) \right\| \\ &+ \left\| Z_{1} \left( \left( P_{U_{2} \perp} \hat{U}_{2}^{(t)} \right) \otimes \left( P_{U_{3} \perp} \hat{U}_{3}^{(t)} \right) \right\| \\ &+ \left\| Z_{1} \left( \left( P_{U_{2} \perp} \hat{U}_{2}^{(t)} \right) \otimes \left( P_{U_{3} \perp} \hat{U}_{3}^{(t)} \right) \right\| \\ &+ \left\| Z_{1} \left( \left( P_{U_{2} \perp} \hat{U}_{2}^{(t)} \right) \otimes \left( P_{U_{3} \perp} \hat{U}_{3}^{(t)} \right) \right\| \\ &+ \left\| Z_{1} \left( \left( P_{U_{2} \perp} \hat{U}_{2}^{(t)} \right) \otimes \left( P_{U_{3} \perp} \hat{U}_{3}^{(t)} \right) \right\| \\ &+ \left\| Z_{1} \left( \left( P_{U_{2} \perp} \hat{U}_{2}^{(t)} \right) \otimes \left( P_{U_{3} \perp} \hat{U}_{3}^{(t)} \right) \right\| \\ &+ \left\| Z_{1} \left( \left( P_{U_{2} \perp} \hat{U}_{2}^{(t)} \right) \otimes \left( P_{U_{3} \perp} \hat{U}_{3}^{(t)} \right) \right\| \\ &+ \left\| Z_{1} \left( \left( P_{U_{2} \perp} \hat{U}_{2}^{(t)} \right) \otimes \left( P_{U_{3} \perp} \hat{U}_{3}^{(t)} \right) \right\| \\ &+ \left\| Z_{1} \left( \left( P_{U_{2} \perp} \hat{U}_{2}^{(t)} \right) \otimes \left( P_{U_{3} \perp} \hat{U}_{3}^{(t)} \right) \right\| \\ &+ \left\| Z_{1} \left( \left( P_{U_{2} \perp} \hat{U}_{2}^{(t)} \right) \otimes \left( P_{U_{3} \perp} \hat{U}_{3}^{(t)$$

 $\leq C_2 \sqrt{p_1} + 3C_1 \sqrt{pr} L_t$  (since the spectral  $\sin \Theta$  norm is at most 1).

Since  $U_1$  and  $\hat{U}_1^{(t+1)}$  are respectively the leading r singular vectors of  $X_1^{(t)}$  and  $Y_1^{(t)}$ , by Wedin's  $\sin \Theta$  theorem [61],

$$\left\| \sin \Theta \left( \hat{U}_{1}^{(t+1)}, U_{1} \right) \right\| \leq \frac{\|Z_{1}^{(t)}\|}{\sigma_{r} \left( X_{1}^{(t)} \right)} \leq \frac{C_{2} \sqrt{p_{1}} + C_{1} \sqrt{p_{r}} L_{t}}{\lambda (1 - L_{t}^{2})}$$

$$\leq \frac{2C_{2} \sqrt{p_{1}}}{\lambda} + \frac{4C_{1} \sqrt{p_{r}}}{\lambda} L_{t}.$$
(32)

We can similarly prove that

$$\left\|\sin\Theta\left(\hat{U}_{2}^{(t+1)}, U_{2}\right)\right\| \leq \frac{2C_{2}\sqrt{p_{2}}}{\lambda} + \frac{4C_{1}\sqrt{pr}}{\lambda}L_{t},$$
$$\left\|\sin\Theta\left(\hat{U}_{3}^{(t+1)}, U_{3}\right)\right\| \leq \frac{2C_{2}\sqrt{p_{3}}}{\lambda} + \frac{4C_{1}\sqrt{pr}}{\lambda}L_{t}.$$

Finally, since  $\max\{p_1, p_2, p_3\} \leq C_0 p$  and  $\max\{r_1, r_2, r_3\} \leq C_0 p^{1/2}$ , there exists a large constant  $C_{gap} > 0$  such that when  $\lambda \geq C_{gap} p^{3/4}$ ,

$$\frac{2C_2\sqrt{p_1}}{\lambda} + \frac{4C_1\sqrt{pr}}{\lambda}L_t \le \frac{1}{2} \quad \text{and} \quad \frac{4C_1\sqrt{pr}}{\lambda} \le \frac{1}{2}$$
 (33)

Then we have finished the proof for (29) for t = 0. By induction, we can sequentially prove that (29) for all  $t \ge 0$ .

At this point, (29) yields

$$L_{t+1} \leq \frac{C\sqrt{p}}{\lambda} + \frac{4C_1\sqrt{pr}}{\lambda}L_t, \qquad t = 1, 2, \dots, t_{\text{max}} - 1$$

$$\Rightarrow L_{t+1} - \frac{2C\sqrt{p}}{\lambda} \leq \frac{4C_1\sqrt{pr}}{\lambda} \left(L_t - \frac{2C\sqrt{p}}{\lambda}\right), \quad \text{(since (33))},$$

$$\Rightarrow L_{t_{\text{max}}} - \frac{2C\sqrt{p}}{\lambda} \leq \left(\frac{4C_1\sqrt{pr}}{\lambda}\right)^{t_{\text{max}}} \cdot \left(L_0 - \frac{2C\sqrt{p}}{\lambda}\right)$$

$$\Rightarrow L_{t_{\text{max}}} \leq \frac{2C\sqrt{p}}{\lambda} + \frac{L_0}{2^{t_{\text{max}}}} = \frac{2C\sqrt{p}}{\lambda} + \frac{1}{2^{t_{\text{max}}}}C\left(\frac{\sqrt{p}\lambda + p^{3/2}}{\lambda^2}\right) \leq \frac{3C\sqrt{p}}{\lambda}$$

when  $t_{\text{max}} \geq C \left( \log \left( \frac{p}{\lambda} \right) \vee 1 \right)$ . Therefore, we have the following upper bound for spectral  $\sin \Theta$  norm loss for  $\hat{U}_k^{t_{\text{max}}} = \hat{U}_k$ ,

$$\|\sin\Theta(\hat{U}_k, U_k)\| \le L_{t_{\max}} \le \frac{C\sqrt{p_k}}{\lambda}.$$
 (34)

when (25), (26), (27) holds.

By the same calculation, we can also prove  $L_{t_{\text{max}}-1}$  satisfies  $L_{t_{\text{max}}-1} \leq C\sqrt{p}/\lambda$ . We prepare this inequality for the use in the next step.

3. In this step, we develop the upper bound for  $\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}}$  under the assumptions of (25), (26), (27), and

$$\left\| \mathbf{Z} \times_{1} \hat{U}_{1}^{\top} \times_{2} \hat{U}_{2}^{\top} \times_{3} \hat{U}_{3}^{\top} \right\|_{F} \leq C \left( \sqrt{p_{1}r_{1}} + \sqrt{p_{2}r_{2}} + \sqrt{p_{3}r_{3}} \right). \tag{35}$$

Instead of working on  $\|\mathbf{X}\|_{\mathrm{F}}$  and  $\hat{U}_{k}^{(t_{\mathrm{max}})}$  directly, we take one step back and work on the evolution of  $\hat{U}_{k}^{(t_{\mathrm{max}}-1)}$  to  $\hat{U}_{k}^{(t_{\mathrm{max}})}$ .

Recall that  $\hat{U}_1 = \hat{U}_1^{(t_{\text{max}})}, \hat{U}_2 = \hat{U}_2^{(t_{\text{max}})}, \hat{U}_3 = \hat{U}_3^{(t_{\text{max}})}; \hat{U}_{1\perp}, \hat{U}_{2\perp}, \hat{U}_{3\perp}$  are the orthogonal complements of  $\hat{U}_1, \hat{U}_2, \hat{U}_3$ , respectively;  $\hat{\mathbf{X}} = \mathbf{Y} \times_1 P_{\hat{U}_1} \times_2 P_{\hat{U}_2} \times_3 P_{\hat{U}_3}$ . In the previous step we have also proved that

$$\left\| \sin \Theta(\hat{U}_k, U_k) \right\| \le C \frac{\sqrt{p_k}}{\lambda}, \quad k = 1, 2, 3.$$

Then we have the following decomposition for the estimation error

$$\begin{aligned} & \left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_{\mathrm{F}} \\ & \leq \left\| \mathbf{X} - \mathbf{X} \times_{1} P_{\hat{U}_{1}} \times_{2} P_{\hat{U}_{2}} \times_{3} P_{\hat{U}_{3}} \right\|_{\mathrm{F}} + \left\| \mathbf{Z} \times_{1} P_{\hat{U}_{1}} \times_{2} P_{\hat{U}_{2}} \times_{3} P_{\hat{U}_{3}} \right\|_{\mathrm{F}} \\ & = \left\| \mathbf{X} \times_{1} P_{\hat{U}_{1\perp}} + \mathbf{X} \times_{1} P_{\hat{U}_{1}} \times_{2} P_{\hat{U}_{2\perp}} + \mathbf{X} \times_{1} P_{\hat{U}_{1}} \times_{2} P_{\hat{U}_{2}} \times_{3} P_{\hat{U}_{3\perp}} \right\|_{\mathrm{F}} \\ & + \left\| \mathbf{Z} \times_{1} \hat{U}_{1}^{\top} \times_{2} \hat{U}_{2}^{\top} \times_{3} \hat{U}_{3}^{\top} \right\|_{\mathrm{F}} \\ & \leq \left\| \mathbf{X} \times_{1} \hat{U}_{1\perp}^{\top} \right\|_{\mathrm{F}} + \left\| \mathbf{X} \times_{2} \hat{U}_{2\perp}^{\top} \right\|_{\mathrm{F}} + \left\| \mathbf{X} \times_{3} \hat{U}_{3\perp}^{\top} \right\|_{\mathrm{F}} \\ & + \left\| \mathbf{Z} \times_{1} \hat{U}_{1}^{\top} \times_{2} \hat{U}_{2}^{\top} \times_{3} \hat{U}_{3}^{\top} \right\|_{\mathrm{F}}. \end{aligned} \tag{36}$$

To obtain the upper bound of  $\|\hat{\mathbf{X}} - \mathbf{X}\|$ , we only need to analyze the four terms in (36) separately. Recall in Step 2, we defined

$$\mathcal{M}_{1}\left(\mathbf{Y} \times_{2} (\hat{U}_{2}^{(t_{\max}-1)})^{\top} \times_{3} (\hat{U}_{3}^{(t_{\max}-1)})^{\top}\right)$$
$$=Y_{1} \cdot \left(\hat{U}_{2}^{(t_{\max}-1)} \otimes \hat{U}_{3}^{(t_{\max}-1)}\right) := Y_{1}^{(t_{\max}-1)},$$

 $X_1^{(t_{\rm max}-1)},\,Z_1^{(t_{\rm max}-1)}$  are defined similarly. Based on the calculation in (30) and (31), we have

$$\sigma_{\min}(X_1^{(t_{\max}-1)}) \ge \sigma_r(X_1) \cdot \sigma_{\min}\left(U_2^{\top} \hat{U}_2^{(t_{\max}-1)}\right) \cdot \sigma_{\min}\left(U_3^{\top} \hat{U}_3^{(t_{\max}-1)}\right) \ge \frac{3}{4}\lambda,$$

$$\|Z_1^{(t_{\max}-1)}\| \le C\sqrt{p_1} + C\sqrt{pr}L_{t_{\max}-1} \le C\sqrt{p_1} + C\sqrt{pr} \cdot \frac{\sqrt{p_1}}{\lambda} \le C\sqrt{p_1}.$$

Since  $\hat{U}_1$  is the leading r left singular vectors of  $Y_1^{(t_{\text{max}}-1)} = X_1^{(t_{\text{max}}-1)} + Z_1^{(t_{\text{max}}-1)}$ , Lemma 6 implies

$$\begin{aligned} & \left\| P_{\hat{U}_{1\perp}} \mathcal{M}_1 \left( \mathbf{X} \times_2 (\hat{U}_2^{(t_{\text{max}} - 1)})^\top \times_3 (\hat{U}_3^{(t_{\text{max}} - 1)})^\top \right) \right\|_{\mathrm{F}} \\ & = \left\| P_{\hat{U}_{1\perp}} X_1^{(t_{\text{max}} - 1)} \right\|_{\mathrm{F}} \le C \sqrt{p_1 r_1}. \end{aligned}$$

As a result,

$$\begin{aligned} & \left\| \mathbf{X} \times_{1} P_{\hat{U}_{1\perp}} \right\|_{F} = \left\| P_{\hat{U}_{1\perp}} \cdot X_{1} \cdot (P_{U_{2}} \otimes P_{U_{3}}) \right\|_{F} = \left\| P_{\hat{U}_{1\perp}} \cdot X_{1} \cdot (U_{2} \otimes U_{3}) \right\|_{F} \\ & \leq \left\| P_{\hat{U}_{1\perp}} X_{1} \left( \hat{U}_{2}^{(t_{\max}-1)} \otimes \hat{U}_{3}^{(t_{\max}-1)} \right) \right\|_{F} \cdot \sigma_{\min}^{-1} (U_{2}^{\top} \hat{U}_{2}^{(t_{\max}-1)}) \cdot \sigma_{\min}^{-1} (U_{3}^{\top} \hat{U}_{3}^{(t_{\max}-1)}) \\ & \leq C \sqrt{p_{1} r_{1}} \frac{1}{\sqrt{1 - (1/2)^{2}}} \frac{1}{\sqrt{1 - (1/2)^{2}}} \leq C \sqrt{p_{1} r_{1}}. \end{aligned}$$
(37)

Similarly, we can show

$$\|\hat{\mathbf{X}} \times_2 P_{\hat{U}_{2\perp}}\|_{\mathbf{F}} \le C\sqrt{p_2 r_2}, \quad \|\hat{\mathbf{X}} \times_3 P_{\hat{U}_{3\perp}}\|_{\mathbf{F}} \le C\sqrt{p_3 r_3}.$$
 (38)

Now combining (36), (35), (37), and (38), we have

$$\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathbf{F}} \le C\left(\sqrt{p_1 r_1} + \sqrt{p_2 r_2} + \sqrt{p_3 r_3}\right)$$
 (39)

for some constant C > 0.

4. We finalize the proof for Theorem 1 in this step. By Lemma 5, we know (26), (27), and (35) hold with probability at least  $1 - C \exp(-cp)$ . By the result in Step 1, we know (25) holds with probability at least  $1 - C \exp(-cp)$ . Let  $Q = \{(26), (27), (35), (25) \text{ all hold}\}$ , then

$$P(Q) \ge 1 - C \exp(-cp). \tag{40}$$

By Steps 2 and 3, one has  $\|\sin\Theta(\hat{U}_k,U_k)\| \leq C\sqrt{p_k}/\lambda, k=1,2,3$ , and

$$\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}} \le C \left(\sqrt{p_1 r_1} + \sqrt{p_2 r_2} + \sqrt{p_3 r_3}\right) \text{ under } Q.$$

It remains to consider situation under  $Q^c$ . By definition,  $\hat{\mathbf{X}}$  is a projection of  $\mathbf{Y}$ , so

$$\|\hat{\mathbf{X}}\|_{F} \le \|\mathbf{Y}\|_{F} \le \|\mathbf{X}\|_{F} + \|\mathbf{Z}\|_{F}.$$

Then we have the following rough upper bound for the 4-th moment of recovery error,

$$\mathbb{E}\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}}^{4} \le C \left(\mathbb{E}\|\hat{\mathbf{X}}\|_{\mathrm{F}}^{4} + \|\mathbf{X}\|_{\mathrm{F}}^{4}\right) \le C\|\mathbf{X}\|_{\mathrm{F}}^{4} + C\mathbb{E}\|\mathbf{Z}\|_{\mathrm{F}}^{4}$$
$$\le C \exp(c_{0}p) + C\mathbb{E}\left(\chi_{p_{1}p_{2}p_{3}}^{2}\right)^{2} \le C \exp(c_{0}p) + Cp^{6}.$$

The the following upper bound holds for the Frobenius norm risk of  $\hat{\mathbf{X}}$ ,

$$\mathbb{E}\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}}^{2} = \mathbb{E}\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}}^{2} \mathbf{1}_{Q} + \mathbb{E}\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}}^{2} \mathbf{1}_{Q^{c}}$$

$$= C\left(\sqrt{p_{1}r_{1}} + \sqrt{p_{2}r_{2}} + \sqrt{p_{3}r_{3}}\right) + \sqrt{\mathbb{E}\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}}^{4} \cdot \mathbb{E}_{Q^{c}}}$$

$$\stackrel{(40)}{\leq} C\left(\sqrt{p_{1}r_{1}} + \sqrt{p_{2}r_{2}} + \sqrt{p_{3}r_{3}}\right) + C\exp\left((c_{0} - c)p/2\right) + Cp^{3}\exp(-cp/2).$$

Thus, one can select  $c_0 < c$  to ensure that

$$\mathbb{E}\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}}^{2} \leq C \left(p_{1}r_{1} + p_{2}r_{2} + p_{3}r_{3}\right).$$

Additionally, since  $\sigma_{r_k}(\mathcal{M}_k(\mathbf{X})) \geq \lambda$ , we have  $\|\mathbf{X}\|_{\mathrm{F}^2} = \|\mathcal{M}_k(\mathbf{X})\|_{\mathrm{F}}^2 \geq r_k \lambda^2$  for k = 1, 2, 3, which implies  $\|\mathbf{X}\|_{\mathrm{F}}^2 \geq \max\{r_1, r_2, r_3\}\lambda = \lambda r$ , then

$$\mathbb{E}\frac{\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}}^2}{\|\mathbf{X}\|_{\mathrm{F}}^2} \le \frac{p_1 + p_2 + p_3}{\lambda}.$$

Moreover, by definition,  $\|\sin\Theta(\hat{U}_k, U_k)\| \le 1$ . Thus we have the following upper bound for the spectral  $\sin\Theta$  risk for  $\hat{U}_k$ ,

$$\mathbb{E}\|\sin\Theta(\hat{U}_{k}, U_{k})\| \leq \mathbb{E}\|\sin\Theta(\hat{U}_{k}, U_{k})\|1_{Q} + \mathbb{E}\|\sin\Theta(\hat{U}_{k}, U_{k})\|1_{Q^{c}}$$

$$= C\frac{\sqrt{p_{k}}}{\lambda} + \sqrt{\mathbb{E}\|\sin\Theta(\hat{U}_{k}, U_{k})\|^{4} \cdot E1_{Q^{c}}} \stackrel{(40)}{\leq} C\frac{\sqrt{p_{k}}}{\lambda} + \sqrt{C\exp(-cp)}.$$

By definition of  $\lambda$ , we know  $\lambda = \sigma_{r_k}(\mathcal{M}_k(\mathbf{X})) \leq \frac{\|\mathbf{X}\|_F}{\sqrt{r_k}} \leq \frac{C \exp(c_0 p)}{\sqrt{r_k}}$ , so one can select small constant  $c_0 > 0$  to ensure that

$$\frac{\sqrt{p_k}}{\lambda} \ge \frac{\sqrt{p_k r_k}}{C \exp(c_0 p)} \ge c \sqrt{\exp(-c p)},$$

which implies  $\mathbb{E}\|\sin\Theta(\hat{U}_k,U_k)\| \leq C\frac{p_k}{\lambda}$ . Finally, we can derive the general Schatten q-sin  $\Theta$ -norm risk via Hölder's inequality,

$$\mathbb{E}r_k^{-1/q}\mathbb{E}\|\sin\Theta(\hat{U}_k,U_k)\|_q \le \mathbb{E}\|\sin\Theta(\hat{U}_k,U_k)\| \le C\frac{\sqrt{p_k}}{\lambda}.$$

Summarizing from Steps 1-4, we have finished the proof of Theorem 1.  $\Box$ 

#### 7.2 Proof of Theorem 4

We particularly show that it suffices to consider sparse tensor models and we set  $\sigma = 1$  for brevity. A tensor  $\mathbf{X} \in \mathbb{R}^{p_1 \times p_2 \times p_3}$  is sparse with respect to parameters  $\mathcal{S}(\mathbf{X}) = (s_1, s_2, s_3)$  if there exists  $S_k(\mathbf{X}) \subset [p_k] := \{1, 2, \dots, p_k\}, k = 1, 2, 3$  such that

$$X_{ijk} = 0, \quad \forall (i, j, k) \in [p_1] \times [p_2] \times [p_3] \setminus S_1(\mathbf{X}) \times S_2(\mathbf{X}) \times S_3(\mathbf{X})$$

with  $|S_k(\mathbf{X})| \leq s_k, k = 1, 2, 3$ . It means that the nonzero entries of  $\mathbf{X}$  are constrained in the block  $S_1(\mathbf{X}) \times S_2(\mathbf{X}) \times S_3(\mathbf{X})$ . Define the subset  $\mathcal{M}(\boldsymbol{p}, k, \boldsymbol{r}, \lambda) \subset \mathcal{F}_{\boldsymbol{p}, \boldsymbol{r}}(\lambda)$  for integer  $k = \lfloor p^{(1-\tau)/2} \rfloor$  as follows,

$$\mathcal{M}(\boldsymbol{p}, k, \boldsymbol{r}, \lambda) := \left\{ \mathbf{X} \in \mathcal{F}_{\boldsymbol{p}, \boldsymbol{r}}(\lambda) : \mathcal{S}(\mathbf{X}) \le (20k, 20k, 20k) \right\},$$

containing sparse tensors in  $\mathcal{F}_{p,r}(\lambda)$ . Consider two disjoint subsets of  $\mathcal{M}(p,k,r,\lambda)$ :

$$\mathcal{M}_0(\boldsymbol{p}, k, \boldsymbol{r}, \lambda) := \Big\{ \mathbf{X} \in \mathcal{M}(\boldsymbol{p}, k, \boldsymbol{r}, \lambda), S_1(\mathbf{X}) \cup S_2(\mathbf{X}) \cup S_3(\mathbf{X}) \subset [p/2] \Big\},\,$$

and

$$\mathcal{M}_1(\boldsymbol{p}, k, \boldsymbol{r}, \lambda) := \left\{ \mathbf{X} \in \mathcal{M}(\boldsymbol{p}, k, \boldsymbol{r}, \lambda), S_1(\mathbf{X}) \cup S_2(\mathbf{X}) \cup S_3(\mathbf{X}) \subset [p] \setminus [p/2] \right\}$$

where matrices in  $\mathcal{M}_0(\boldsymbol{p}, k, \boldsymbol{r}, \lambda)$  and  $\mathcal{M}_1(\boldsymbol{p}, k, \boldsymbol{r}, \lambda)$  are supported on disjoint blocks, so are their singular vectors. Given the observation:

$$\mathbf{Y} = \mathbf{X} + \mathbf{Z} \in \mathbb{R}^{p_1 \times p_2 \times p_3}.$$

the following testing problem is studied:

$$H_0: \mathbf{X} \in \mathcal{M}_0(\boldsymbol{p}, k, \boldsymbol{r}, \lambda) \quad \text{V.S.} \quad H_1: \mathbf{X} \in \mathcal{M}_1(\boldsymbol{p}, k, \boldsymbol{r}, \lambda).$$
 (41)

A test is then defined as  $\phi(\cdot): \mathbb{R}^{p_1 \times p_2 \times p_3} \to \{0,1\}$  whose risk is given as

$$\mathcal{R}_{\boldsymbol{p},\boldsymbol{r},\lambda}(\phi) = \sup_{\mathbf{X} \in \mathcal{M}_0(\boldsymbol{p},k,\boldsymbol{r},\lambda)} \mathbb{P}_{\mathbf{X}} \{ \phi(\mathbf{Y}) = 1 \} + \sup_{\mathbf{X} \in \mathcal{M}_1(\boldsymbol{p},k,\boldsymbol{r},\lambda)} \mathbb{P}_{\mathbf{X}} \{ \phi(\mathbf{Y}) = 0 \},$$

the worst case of Type-I+II error.

**Lemma 1.** Suppose Hypothesis  $\mathbf{H}(\tau)$  for some  $\tau \in (0,1)$ . Let  $\{\phi_p\}$  be any sequence of polynomial-time tests of (41). There exists an absolute constant  $c_0 > 0$  such that if  $\lambda \leq c_0 \left(\frac{p^{3(1-\tau)/4}}{\sqrt{\log p}}\right)$ , then as long as  $\min\{r_1, r_2, r_3\} \geq 1$ ,

$$\liminf_{p \to \infty} \mathcal{R}_{\boldsymbol{p}, \boldsymbol{r}, \lambda}(\phi_p) \ge \frac{1}{2}.$$

Now we move back to the proof for Theorem 4. Suppose that, on the contradiction, for any k=1,2,3, there exists a sub-sequence  $(\hat{U}_k^{(p)})$  such that

$$\lim_{p \to \infty} \sup_{\mathbf{X} \in \mathcal{F}_{p,r}(\lambda)} \mathbb{E} \| \sin \Theta (\hat{U}_k^{(p)}, U_k(\mathbf{X})) \| = 0,$$

which implies that

$$\lim_{p \to \infty} \sup_{\mathbf{X} \in \mathcal{F}_{\mathbf{p}, \mathbf{r}}(\lambda)} \mathbb{P}\left( \left\| U_k U_k^\top - \hat{U}_k^{(p)} \left( \hat{U}_k^{(p)} \right)^\top \right\| \le \frac{1}{3} \right) = 1.$$
 (42)

Define a sequence of tests  $\phi_p : \mathbb{R}^{p_1 \times p_2 \times p_3} \mapsto \{0,1\}$  as

$$\phi_p(\mathbf{Y}) := \begin{cases} 0, & \text{if } \left\| (\hat{U}_k^{(p)})_{[1:p/2,:]} (\hat{U}_k^{(p)})_{[1:p/2,:]}^\top \right\| \ge \frac{2}{3}, \\ 1, & \text{otherwise,} \end{cases}$$

where  $(\hat{U}_k^{(p)})_{[1:p/2,:]}$  denote the first p/2 rows of  $\hat{U}_k^{(p)}$ . Clearly,

$$\mathcal{R}_{\boldsymbol{p},\boldsymbol{r},\lambda}(\phi_p) \leq \sup_{\mathbf{X} \in \mathcal{M}_0(\boldsymbol{p},k,\boldsymbol{r},\lambda)} \mathbb{P}_{\mathbf{X}} \Big( \| U_k U_k^\top - \hat{U}_k^{(p)} (\hat{U}_k^{(p)})^\top \| > \frac{1}{3} \Big)$$

$$+ \sup_{\mathbf{X} \in \mathcal{M}_1(\boldsymbol{p},k,\boldsymbol{r},\lambda)} \mathbb{P}_{\mathbf{X}} \Big( \| U_k U_k^\top - \hat{U}_k^{(p)} \big( \hat{U}_k^{(p)} \big)^\top \| \ge \frac{2}{3} \Big),$$

which implies  $\lim_{p\to\infty} \mathcal{R}_{\boldsymbol{p},\boldsymbol{r},\lambda}(\phi_p) = 0$ , contradicting Lemma 1. Now, we prove claim (18). Suppose that, on the contradiction, there exists a sub-sequence  $(\hat{\mathbf{X}}^{(p)})$  such that

$$\lim_{p \to \infty} \sup_{\mathbf{X} \in \mathcal{F}_{\mathbf{p}, \mathbf{r}}(\lambda)} \mathbb{E} \frac{\|\hat{\mathbf{X}}^{(p)} - \mathbf{X}\|_{\mathrm{F}}^2}{\|\mathbf{X}\|_{\mathrm{F}}^2} = 0,$$

which implies

$$\lim_{p \to \infty} \sup_{\mathbf{X} \in \mathcal{F}_{p,r}(\lambda)} \mathbb{P}\left( \|\hat{\mathbf{X}}^{(p)} - \mathbf{X}\|_{F} \le \frac{1}{3} \|\mathbf{X}\|_{F} \right) = 1.$$

$$(43)$$

Define a sequence of test  $\phi_p : \mathbb{R}^{p_1 \times p_2 \times p_3} \mapsto \{0,1\}$  as

$$\phi_p(\mathbf{Y}) := \begin{cases} 0, & \text{if } \|(\hat{\mathbf{X}}^{(p)})_{[V_1, V_1, V_1]}\|_{\mathcal{F}} \ge \|(\hat{\mathbf{X}}^{(p)})_{[V_2, V_2, V_2]}\|_{\mathcal{F}}, \\ 1, & \text{otherwise,} \end{cases}$$

where  $V_1 = [p/2], V_2 = [p] \setminus V_1$  and  $\mathbf{X}_{[V_1, V_1, V_1]}$  denotes the sub-tensor on the block  $V_1 \times V_1 \times V_1$ . Under  $H_0$ , if  $\phi_p(\mathbf{Y}) = 1$ , then

$$\begin{split} \|\hat{\mathbf{X}}^{(p)} - \mathbf{X}\|_{\mathrm{F}}^2 &\geq \|\left(\hat{\mathbf{X}}^{(p)} - \mathbf{X}\right)_{[V_1, V_1, V_1]}\|_{\mathrm{F}}^2 + \|\left(\hat{\mathbf{X}}^{(p)}\right)_{[V_2, V_2, V_2]}\|_{\mathrm{F}}^2 \\ &\geq \left\|\left(\hat{\mathbf{X}}^{(p)} - \mathbf{X}\right)_{[V_1, V_1, V_1]}\right\|_{\mathrm{F}}^2 + \left\|\left(\hat{\mathbf{X}}^{(p)}\right)_{[V_1, V_1, V_1]}\right\|_{\mathrm{F}}^2 &\geq \frac{1}{2}\|\mathbf{X}\|_{\mathrm{F}}^2. \end{split}$$

Clearly,

$$\mathcal{R}_{\boldsymbol{p},\boldsymbol{r},\lambda}(\phi_p) = \sup_{\mathbf{X} \in \mathcal{M}_0(\boldsymbol{p},k,\boldsymbol{r},\lambda)} \mathbb{P}\left(\|\hat{\mathbf{X}}^{(p)} - \mathbf{X}\|_F \ge \frac{\sqrt{2}}{2} \|\mathbf{X}\|_F\right) + \sup_{\mathbf{X} \in \mathcal{M}_1(\boldsymbol{p},k,\boldsymbol{r},\lambda)} \mathbb{P}\left(\|\hat{\mathbf{X}}^{(p)} - \mathbf{X}\|_F \ge \frac{\sqrt{2}}{2} \|\mathbf{X}\|_F\right),$$

which implies that  $\mathcal{R}_{p,r,\lambda}(\phi_p) \to 0$  as  $p \to \infty$  based on (43), which contradicts Lemma 1.  $\square$ 

#### 7.3 Proof of Theorems 3

Without loss of generality we can assume  $\sigma = 1$  throughout the proof. First, we construct the core tensor  $\tilde{\mathbf{S}} \in \mathbb{R}^{r_1 \times r_2 \times r_3}$  with i.i.d. standard Gaussian entries, then according to random matrix theory (c.f. Corollary 5.35 in [62]), with probability at least  $1 - 6e^{-x}$ , we have

$$\sqrt{r_{k+1}r_{k+2}} - \sqrt{r_k} - x \le \sigma_{\min}(\mathcal{M}_k(\tilde{\mathbf{S}})) \le \sigma_{\max}(\mathcal{M}_k(\tilde{\mathbf{S}})) \le \sqrt{r_{k+1}r_{k+2}} + \sqrt{r_k} + x,$$

for k = 1, 2, 3. Plug in x = 1.8, by simple calculation, we can see there is a positive probability that

$$\sqrt{r_{k+1}r_{k+2}} - \sqrt{r_k} - 1.8 \le \sigma_{\min}(\mathcal{M}_k(\tilde{\mathbf{S}})) \le \sigma_{\max}(\mathcal{M}_k(\tilde{\mathbf{S}})) \le \sqrt{r_{k+1}r_{k+2}} + \sqrt{r_k} + 1.8. \tag{44}$$

Note that

$$\begin{split} r_1 r_2 & \geq 4 r_3, r_2 r_3 \geq 4 r_1, r_3 r_1 \geq 4 r_2, \quad \Rightarrow \quad r_1 r_2 r_3 \geq 4 \max_{1 \leq k \leq 3} \{r_k\}^2 \\ & \Rightarrow \quad r_k \frac{r_{k+1}}{\max\{r_k\}} \frac{r_{k+2}}{\max\{r_k\}} \geq 4, \Rightarrow r_k \geq 4 \\ & \Rightarrow \quad r_{k+1} r_{k+2} \geq \frac{4 \max\{r_k\}^2}{r_k} \geq 4 r_k \geq 16, \end{split}$$

we know

$$\begin{split} &\sqrt{r_{k+1}r_{k+2}}-\sqrt{r_k}-1.8\\ \ge &\sqrt{r_{k+1}r_{k+2}}-\sqrt{\frac{r_{k+1}r_{k_2}}{4}}-\frac{1.8}{4}\sqrt{r_{k+1}r_{k+2}}=0.05\sqrt{r_{k+1}r_{k+2}};\\ &\sqrt{r_{k+1}r_{k+2}}-\sqrt{r_k}+1.8\\ \le &\sqrt{r_{k+1}r_{k+2}}+\sqrt{\frac{r_{k+1}r_{k_2}}{4}}+\frac{1.8}{4}\sqrt{r_{k+1}r_{k+2}}=1.95\sqrt{r_{k+1}r_{k+2}}. \end{split}$$

By previous arguments, there exists  $\tilde{\mathbf{S}} \in \mathbb{R}^{r_1 \times r_2 \times r_3}$  such that  $c\sqrt{r_{k+1}r_{k+2}} \leq \sigma_{\min}(\mathcal{M}_k(\tilde{\mathbf{S}})) \leq C\sqrt{r_{k+1}r_{k+2}}$  for k = 1, 2, 3. Now, we construct the scaled core tensor  $\mathbf{S} = \tilde{\mathbf{S}} \frac{\lambda}{\min_{k=1,2,3} \sigma_{\min}(\mathcal{M}_k(\tilde{\mathbf{S}}))}$ . Given  $r \leq r_1, r_2, r_3 \leq C_0 r$ , we know  $\mathbf{S} \in \mathbb{R}^{r_1 \times r_2 \times r_3}$  satisfies the following property

$$\lambda \le \sigma_{\min}(\mathcal{M}_k(\mathbf{S})) \le \sigma_{\max}(\mathcal{M}_k(\mathbf{S})) \le C\lambda, \quad k = 1, 2, 3.$$
 (45)

Proof of the first claim. It suffices to consider k = 1. We construct a large subset of  $\mathbb{O}_{p_1,r_1}$  whose elements are well separated in Schatten q-norms for all  $1 \leq q \leq +\infty$ . To this end, we need some preliminary facts about the packing number in Grassmann manifold  $\mathcal{G}_{p,r}$ , which is the set of all r-dimensional subspaces of  $\mathbb{R}^p$ . Given such a subspace  $L \subset \mathbb{R}^p$  with  $\dim(L) = r$ , let  $U_L \in \mathbb{O}(p,r)$  denote the orthonormal basis of L. Denote  $\mathcal{B}_{p,r} := \{U_L, L \in \mathcal{G}_{p,r}\}$  which is actually a subset of  $\mathbb{O}_{p,r}$  and will be equipped with Schatten q-norm distances for all  $q \in [1, +\infty]$ :  $d_q(U_{L_1}, U_{L_2}) := \|U_{L_1}U_{L_1}^\top - U_{L_2}U_{L_2}^\top\|_q$ . Recall that the  $\varepsilon$ -packing number of a metric space (T, d) is defined as

$$D(T,d,\varepsilon) := \max \Big\{ n : \text{there are } t_1,\ldots,t_n \in T, \text{ such that } \min_{\mathbf{i} \neq \mathbf{j}} \ \mathrm{d}(\mathbf{t_i},\mathbf{t_j}) > \varepsilon \Big\}.$$

The following lemma can be found in Lemma 5 in [63] which controls the packing numbers of  $\mathcal{B}_{p,r}$  with respect to Schatten distances  $d_q$ .

**Lemma 2.** For all integers  $1 \le r \le p$  such that  $r \le p - r$ , and all  $1 \le q \le +\infty$ , the following bound holds

$$\left(\frac{c}{\varepsilon}\right)^{r(p-r)} \le D\left(\mathcal{B}_{p,r}, d_q, \varepsilon r^{1/q}\right) \le \left(\frac{C}{\varepsilon}\right)^{r(p-r)}$$

with absolute constants c, C > 0.

We are in position to construct a well-separated subset of  $\mathbb{O}_{p_1,r_1}$ . According to Lemma 2 by choosing  $\varepsilon = \frac{c}{2}$ , there exists a subset  $\mathcal{V}_{p_1-r_1,r_1} \subset \mathbb{O}_{p_1-r_1,r_1}$  with  $\operatorname{Card}(\mathcal{V}_{p_1-r_1,r_1}) \geq 2^{r_1(p_1-2r_1)}$  such that for each  $V_1 \neq V_2 \in \mathcal{V}_{p_1-r_1,r_1}$ ,

$$||V_1V_1^{\top} - V_2V_2^{\top}||_q \ge \frac{c}{2}r_1^{1/q}.$$

Now, fix a  $\delta > 0$  whose value is to be determined later. For every  $V \in \mathcal{V}_{p_1-r_1,r_1}$ , define  $\tilde{V} \in \mathbb{O}_{p_1,r_1}$  as follows

$$\tilde{V} = \left(\begin{array}{c} \sqrt{1 - \delta} I_{r_1} \\ \sqrt{\delta} V \end{array}\right).$$

It is easy to check that  $\tilde{V} \in \mathbb{O}_{p_1,r_1}$  as long as  $V \in \mathbb{O}_{p_1-r_1,r_1}$ . We conclude with a subset  $\mathcal{V}_{p_1,r_1} \subset \mathbb{O}_{p_1,r_1}$  with  $\operatorname{Card}(\mathcal{V}_{p_1,r_1}) = \operatorname{Card}(\mathcal{V}_{p_1-r_1,r_1}) \geq 2^{r_1(p_1-2r_1)}$ . Moreover, for  $\tilde{V}_1 \neq \tilde{V}_2 \in \mathcal{V}_{p_1,r_1}$ ,

$$\|\tilde{V}_1\tilde{V}_1^\top - \tilde{V}_2\tilde{V}_2^\top\|_q \ge \sqrt{\delta(1-\delta)}\|V_1 - V_2\|_q \ge \frac{c}{2}\sqrt{\delta(1-\delta)}r_1^{1/q}.$$

Meanwhile,

$$\|\tilde{V}_1 - \tilde{V}_2\|_{\mathrm{F}} < \sqrt{\delta} \|V_1 - V_2\|_{\mathrm{F}} < \sqrt{2\delta r_1}.$$

Then we construct a series of fixed signal tensors:  $\mathbf{X}_i = \mathbf{S} \times_1 \tilde{V}_i \times_2 U_2 \times_3 U_3, i = 1, \dots, m$ , where  $U_2 \in \mathbb{O}_{p_2,r_2}.U_3 \in \mathbb{O}_{p_3,r_3}$  are any fixed orthonormal columns,  $\tilde{V}_i \in \mathcal{V}_{p_1,r_1} \subset \mathbb{O}_{p_1,r_1}$  and  $m = 2^{r_1(p_1-2r_1)}$ . By such construction,  $\sigma^2_{\min}(\mathcal{M}_k(\mathbf{X}_i)) \geq \lambda$  for k = 1, 2, 3, so that  $\{\mathbf{X}_i\}_{i=1}^m \subseteq \mathcal{F}_{p,r}(\lambda)$ .

We further let  $\mathbf{Y}_i = \mathbf{X}_i + \mathbf{Z}_i$ , where  $\mathbf{Z}_i$  are i.i.d. standard normal distributed tensors, which implies  $\mathbf{Y}_i \sim N(\mathbf{X}_i, I_{p_1 \times p_2 \times p_3})$ . Then the Kullback-Leibler divergence between the distribution  $\mathbf{Y}_i$  and  $\mathbf{Y}_j$  is

$$D_{\text{KL}}(\mathbf{Y}_{i}||\mathbf{Y}_{j}) = \frac{1}{2} \|\mathbf{X}_{i} - \mathbf{X}_{j}\|_{\text{F}}^{2} = \frac{1}{2} \|\mathbf{S} \times_{1} (\tilde{V}_{i} - \tilde{V}_{j}) \times_{2} U_{2} \times_{3} U_{3}\|_{\text{F}}^{2}$$
$$= \frac{1}{2} \|(\tilde{V}_{i} - \tilde{V}_{j}) \cdot \mathcal{M}_{1}(\mathbf{S}) \cdot (U_{2} \otimes U_{3})^{\top}\|_{\text{F}}^{2} \leq C\lambda^{2} \|\tilde{V}_{i} - \tilde{V}_{j}\|_{\text{F}}^{2} \leq 2C\lambda^{2}\delta r_{1}.$$

Then the generalized Fano's lemma yields the following lower bound

$$\inf_{\hat{U}_1} \sup_{U_1 \in \{V_i\}_{i=1}^m} \mathbb{E} \left\| \hat{U}_1 \hat{U}_1^\top - U_1 U_1^\top \right\|_q \ge \frac{c}{2} \sqrt{\delta(1-\delta)} r_1^{1/q} \left( 1 - \frac{C\lambda^2 \delta r_1 + \log 2}{r_1(p_1 - 2r_1) \log 2} \right). \tag{46}$$

By setting  $\delta = c_1 \frac{(p_1 - 2r_1)}{\lambda^2}$  for a small but absolute constant  $c_1 > 0$ , we obtain

$$\frac{c}{2}\sqrt{\delta(1-\delta)}r_1^{1/q}\left(1-\frac{C\lambda^2\delta r_1+\log 2}{r_1(p_1-2r_1)\log 2}\right) \ge c_0\frac{(p_1-2r_1)^{1/2}}{\lambda}r_1^{1/q}$$

for an absolute constant  $c_0 > 0$ . Then, if  $p_1 \ge 3r_1$ ,

$$\inf_{\hat{U}_1} \sup_{U_1 \in \{V_i\}_{i=1}^m} \mathbb{E} \left\| \hat{U}_1 \hat{U}_1^\top - U_1 U_1^\top \right\|_q \ge c_0 \left( \frac{\sqrt{p_1}}{\lambda} r_1^{1/q} \wedge r_1^{1/q} \right)$$

where  $r_1^{1/q}$  is a trivial term. The first claim in Theorem 3 it thus obtained by viewing the equivalence between the Schatten q-norms and  $\sin\Theta$  Schatten q-norms, see Lemma 3 in the Appendix.

**Proof of second and third claims.** To prove the minimax lower bounds in estimating  $\mathbf{X}$ , we need a different construction scheme. Specifically, we consider the metric space  $(\mathbb{O}_{p_1,r_1}, \|\sin\Theta(\cdot,\cdot)\|_2)$ , fix an  $V_0 \in \mathbb{O}_{p_1,r_1}$ , and consider the following ball of radius  $\varepsilon > 0$  and center  $V_0$ :

$$B(V_0, \varepsilon) = \left\{ V' \in \mathbb{O}_{p_1, r_1} : \| \sin \Theta(V', V) \|_2 \le \varepsilon \right\}.$$

By Lemma 1 in [10], for  $0 < \alpha < 1$  and  $0 < \varepsilon \le 1$ , there exists  $V_1', \dots, V_m' \subseteq B(V_0, \varepsilon)$  such that

$$m \ge \left(\frac{c_0}{\alpha}\right)^{r_1(p_1-r_1)}, \quad \min_{1 \le i < j \le m} \left\|\sin\Theta(V_i', V_j')\right\|_2 \ge \alpha \varepsilon.$$

By the property of  $\sin \Theta$  distance (Lemma 1 in [26]), we can find a rotation matrix  $O_i \in \mathbb{O}_{r_1}$  such that

$$||V_0 - V_i' O_i||_{\mathcal{F}} \le \sqrt{2} ||\sin \Theta(V_0, V_i')||_2 \le \sqrt{2\varepsilon}.$$

We denote  $V_i = V_i'O_i$ , then

$$||V_i - V_0||_{\mathcal{F}} \le \sqrt{2\varepsilon}, \quad ||\sin\Theta(V_i, V_j)||_2 \ge \alpha\varepsilon, \quad 1 \le i < j \le m. \tag{47}$$

Construct  $\mathbf{X}_i = \mathbf{S} \times_1 V_i \times_2 U_2 \times U_3$  for i = 1, ..., m in a similar fashion as above. Then the class of low-rank tensors satisfy the following properties,

$$\|\mathbf{X}_{i} - \mathbf{X}_{j}\|_{F}^{2} = \|\mathbf{S} \times_{1} (V_{i} - V_{j}) \times_{2} U_{2} \times_{3} U_{3}\|_{F}^{2}$$

$$= \frac{1}{2} \|(V_{i} - V_{j}) \cdot \mathcal{M}_{1}(\mathbf{S}) \cdot (U_{2} \otimes U_{3})^{\top}\|_{F}^{2} \geq \frac{1}{2} \sigma_{r_{1}} (\mathcal{M}_{1}(\mathbf{S})) \|V_{i} - V_{j}\|_{F}^{2}$$

$$\geq \frac{\lambda^{2}}{2} \min_{O \in \mathbb{O}_{r}} \|V_{i} - V_{j}O\|_{F}^{2} \quad \text{(by Lemma 1 in [26])}$$

$$\geq \frac{\lambda^{2}}{2} \|\sin \Theta (V_{i}, V_{j})\|_{2}^{2} \geq \alpha^{2} \varepsilon^{2} \lambda^{2}, \quad 1 \leq i < j \leq m.$$

$$(48)$$

$$\min_{1 \le i \le m} \|\mathbf{X}_i\|_{F} = \|\mathbf{S} \times_1 V_i \times_2 U_2 \times_3 U_3\|_{F} = \|\mathbf{S}\|_{F} \ge \|\mathcal{M}_1(\mathbf{S})\|_{F} \ge \lambda \sqrt{r_1}.$$
 (49)

Moreover, under the same model  $\mathbf{Y}_i = \mathbf{X}_i + \mathbf{Z}_i$  as above, the KL-divergence between the distributions of  $\mathbf{Y}_i$  and  $\mathbf{Y}_j$  is

$$D_{\text{KL}}(\mathbf{Y}_{i}||\mathbf{Y}_{j}) = \frac{1}{2} \|\mathbf{X}_{i} - \mathbf{X}_{j}\|_{\text{F}}^{2} = \frac{1}{2} \|\mathbf{S} \times_{1} (V_{i} - V_{j}) \times_{2} U_{2} \times_{3} U_{3}\|_{\text{F}}^{2}$$
$$= \frac{1}{2} \|(V_{i} - V_{j}) \cdot \mathcal{M}_{1}(\mathbf{S}) \cdot (U_{2} \otimes U_{3})^{\top}\|_{\text{F}}^{2} \leq C\lambda^{2} \|V_{i} - V_{j}\|_{\text{F}}^{2} \leq C\lambda^{2} \varepsilon^{2}.$$

Then the generalized Fano's lemma yields the following lower bound

$$\inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X} \in \{\mathbf{X}_i\}_{i=1}^m} \mathbb{E} \left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_{\mathrm{F}}^2 \ge \lambda \alpha \varepsilon \left( 1 - \frac{C \lambda^2 \varepsilon^2 + \log 2}{r_1 (p_1 - r_1) \log(c_0 / \alpha)} \right). \tag{50}$$

By setting  $\varepsilon = \sqrt{\frac{r_1(p_1-r_1)}{2C\lambda^2}} \wedge \sqrt{2r_1}$ ,  $\alpha = (c_0 \wedge 1)/8$ , we have

$$\alpha \varepsilon \left( 1 - \frac{C\lambda^2 \varepsilon^2 + \log 2}{r_1(p_1 - r_1)\log(c_0/\alpha)} \right) \ge c_1 \left( \frac{\sqrt{r_1 p_1}}{\lambda} \wedge \sqrt{r_1} \right)$$

for some small constant  $c_1 > 0$ . Moreover,

$$\inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X} \in \mathcal{F}_{\boldsymbol{p}, \boldsymbol{r}}(\lambda)} \mathbb{E} \left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_{\mathrm{F}}^{2} \ge \inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X} \in \mathcal{F}_{\boldsymbol{p}, \boldsymbol{r}}(\lambda \vee \sqrt{p_{1}})} \mathbb{E} \left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_{\mathrm{F}}^{2}$$

$$\ge \inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X} \in \{\mathbf{X}_{1}, \dots, \mathbf{X}_{m}\}} \mathbb{E} \left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_{\mathrm{F}}^{2} \ge c_{1}(\lambda^{2} \vee p_{1}) \left( \frac{r_{1}p_{1}}{\lambda^{2} \vee p_{1}} \wedge r_{1} \right) \ge c_{1}p_{1}r_{1}.$$
(51)

$$\inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X} \in \mathcal{F}_{\boldsymbol{p}, \boldsymbol{r}}(\lambda)} \mathbb{E} \frac{\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}}^{2}}{\|\mathbf{X}\|_{\mathrm{F}}^{2}} \ge \inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X} \in \mathcal{F}_{\boldsymbol{p}, \boldsymbol{r}}(\lambda \vee p_{1})} \mathbb{E} \frac{\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}}^{2}}{\|\mathbf{X}\|_{\mathrm{F}}^{2}}$$

$$\ge \inf_{\hat{\mathbf{X}}} \sup_{\mathbf{X} \in \{\mathbf{X}_{1}, \dots, \mathbf{X}_{m}\}} \mathbb{E} \frac{\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}}^{2}}{\max_{1 \le i \le m} \|\mathbf{X}\|_{\mathrm{F}}^{2}} \stackrel{(49)(50)}{\ge} \frac{c_{1} \lambda^{2}}{\lambda^{2} r_{1}} \left(\frac{r_{1} p_{1}}{\lambda^{2}} \wedge r_{1}\right) \ge c_{1} \left(\frac{p_{1}}{\lambda^{2}} \wedge 1\right). \tag{52}$$

Finally, we apply the same argument of (50), (51), and (52) to  $U_2, U_3$ , then we can obtain (13) and (14).  $\square$ 

# Acknowledgment

The authors would like to thank Ming Yuan and Zongming Ma for helpful discussions. The authors would also like to thank Genevera Allen and Nathaniel Helwig for pointing out several key references. The authors would also like to thank editors and anomalous referees for your helpful comments and suggestions on improving the paper.

## References

- [1] H. Shen and J. Z. Huang, "Sparse principal component analysis via regularized low rank matrix approximation," *Journal of multivariate analysis*, vol. 99, no. 6, pp. 1015–1034, 2008.
- [2] M. Lee, H. Shen, J. Z. Huang, and J. Marron, "Biclustering via sparse singular value decomposition," *Biometrics*, vol. 66, no. 4, pp. 1087–1095, 2010.
- [3] D. Yang, Z. Ma, and A. Buja, "A sparse singular value decomposition method for high-dimensional data," *Journal of Computational and Graphical Statistics*, vol. 23, no. 4, pp. 923–942, 2014.
- [4] D. Yang, Z. Ma, and A. Buja, "Rate optimal denoising of simultaneously sparse and low rank matrices," *Journal of Machine Learning Research*, vol. 17, no. 92, pp. 1–27, 2016.
- [5] E. J. Candes, C. A. Sing-Long, and J. D. Trzasko, "Unbiased risk estimates for singular value thresholding and spectral estimators," *IEEE transactions on signal processing*, vol. 61, no. 19, pp. 4643–4657, 2013.
- [6] A. A. Shabalin and A. B. Nobel, "Reconstruction of a low-rank matrix in the presence of gaussian noise," *Journal of Multivariate Analysis*, vol. 118, pp. 67–76, 2013.
- [7] D. Donoho, M. Gavish, et al., "Minimax risk of matrix denoising by singular value thresholding," The Annals of Statistics, vol. 42, no. 6, pp. 2413–2440, 2014.
- [8] M. Gavish and D. L. Donoho, "The optimal hard threshold for singular values is  $4/\sqrt{3}$ ," *IEEE Transactions on Information Theory*, vol. 60, no. 8, pp. 5040–5053, 2014.

- [9] H. Zou, T. Hastie, and R. Tibshirani, "Sparse principal component analysis," *Journal of computational and graphical statistics*, vol. 15, no. 2, pp. 265–286, 2006.
- [10] T. T. Cai, Z. Ma, Y. Wu, et al., "Sparse pca: Optimal rates and adaptive estimation," The Annals of Statistics, vol. 41, no. 6, pp. 3074–3110, 2013.
- [11] A. Birnbaum, I. M. Johnstone, B. Nadler, and D. Paul, "Minimax bounds for sparse pca with noisy high-dimensional data," *Annals of statistics*, vol. 41, no. 3, p. 1055, 2013.
- [12] E. J. Candès, X. Li, Y. Ma, and J. Wright, "Robust principal component analysis?," *Journal of the ACM (JACM)*, vol. 58, no. 3, p. 11, 2011.
- [13] H. Zhou, L. Li, and H. Zhu, "Tensor regression with applications in neuroimaging data analysis," *Journal of the American Statistical Association*, vol. 108, no. 502, pp. 540–552, 2013.
- [14] A. Zhang, "Cross: Efficient low-rank tensor completion," arXiv preprint arXiv:1611.01129, 2016.
- [15] A. Karatzoglou, X. Amatriain, L. Baltrunas, and N. Oliver, "Multiverse recommendation: n-dimensional tensor factorization for context-aware collaborative filtering," in *Proceedings* of the fourth ACM conference on Recommender systems, pp. 79–86, ACM, 2010.
- [16] S. Rendle and L. Schmidt-Thieme, "Pairwise interaction tensor factorization for personalized tag recommendation," in *Proceedings of the third ACM international conference on Web search and data mining*, pp. 81–90, ACM, 2010.
- [17] H. Lu, K. N. Plataniotis, and A. N. Venetsanopoulos, "Mpca: Multilinear principal component analysis of tensor objects," *IEEE Transactions on Neural Networks*, vol. 19, no. 1, pp. 18–39, 2008.
- [18] J. Liu, P. Musialski, P. Wonka, and J. Ye, "Tensor completion for estimating missing values in visual data," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 1, pp. 208–220, 2013.

- [19] A. Anandkumar, R. Ge, D. J. Hsu, S. M. Kakade, and M. Telgarsky, "Tensor decompositions for learning latent variable models.," *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 2773–2832, 2014.
- [20] A. Anandkumar, R. Ge, D. Hsu, and S. M. Kakade, "A tensor approach to learning mixed membership community models," *J. Mach. Learn. Res.*, vol. 15, pp. 2239–2312, Jan. 2014.
- [21] N. Li and B. Li, "Tensor completion for on-board compression of hyperspectral images," in 2010 IEEE International Conference on Image Processing, pp. 517–520, IEEE, 2010.
- [22] T. Liu, M. Yuan, and H. Zhao, "Characterizing spatiotemporal transcriptome of human brain via low rank tensor decomposition," arXiv preprint arXiv:1702.07449, 2017.
- [23] E. Richard and A. Montanari, "A statistical model for tensor pca," in *Advances in Neural Information Processing Systems*, pp. 2897–2905, 2014.
- [24] S. B. Hopkins, J. Shi, and D. Steurer, "Tensor principal component analysis via sum-of-square proofs.," in *COLT*, pp. 956–1006, 2015.
- [25] A. Perry, A. S. Wein, and A. S. Bandeira, "Statistical limits of spiked tensor models," arXiv preprint arXiv:1612.07728, 2016.
- [26] T. T. Cai and A. Zhang, "Rate-optimal perturbation bounds for singular subspaces with applications to high-dimensional statistics," *The Annals of Statistics*, vol. to appear, 2017.
- [27] C. J. Hillar and L.-H. Lim, "Most tensor problems are np-hard," *Journal of the ACM* (*JACM*), vol. 60, no. 6, p. 45, 2013.
- [28] L. De Lathauwer, B. De Moor, and J. Vandewalle, "A multilinear singular value decomposition," SIAM journal on Matrix Analysis and Applications, vol. 21, no. 4, pp. 1253–1278, 2000.
- [29] L. De Lathauwer, B. De Moor, and J. Vandewalle, "On the best rank-1 and rank-(r 1, r 2,..., rn) approximation of higher-order tensors," SIAM Journal on Matrix Analysis and Applications, vol. 21, no. 4, pp. 1324–1342, 2000.

- [30] M. A. O. Vasilescu and D. Terzopoulos, "Multilinear subspace analysis of image ensembles," in Computer Vision and Pattern Recognition, 2003. Proceedings. 2003 IEEE Computer Society Conference on, vol. 2, pp. II–93, IEEE, 2003.
- [31] B. N. Sheehan and Y. Saad, "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*, pp. 355–365, SIAM, 2007.
- [32] R. Costantini, L. Sbaiz, and S. Susstrunk, "Higher order svd analysis for dynamic texture synthesis," *IEEE Transactions on Image Processing*, vol. 17, no. 1, pp. 42–52, 2008.
- [33] M. Haardt, F. Roemer, and G. Del Galdo, "Higher-order svd-based subspace estimation to improve the parameter estimation accuracy in multidimensional harmonic retrieval problems," *IEEE Transactions on Signal Processing*, vol. 56, no. 7, pp. 3198–3213, 2008.
- [34] Y. Liu, F. Shang, W. Fan, J. Cheng, and H. Cheng, "Generalized higher-order orthogonal iteration for tensor decomposition and completion," in *Advances in Neural Information Processing Systems*, pp. 1763–1771, 2014.
- [35] Q. Zheng and R. Tomioka, "Interpolating convex and non-convex tensor decompositions via the subspace norm," in Advances in Neural Information Processing Systems, pp. 3106–3113, 2015.
- [36] A. Anandkumar, Y. Deng, R. Ge, and H. Mobahi, "Homotopy analysis for tensor pca," arXiv preprint arXiv:1610.09322, 2016.
- [37] A. Anandkumar, R. Ge, and M. Janzamin, "Guaranteed non-orthogonal tensor decomposition via alternating rank-1 updates," arXiv preprint arXiv:1402.5180, 2014.
- [38] W. W. Sun, J. Lu, H. Liu, and G. Cheng, "Provable sparse tensor decomposition," *Journal of Royal Statistical Association*, 2015.
- [39] G. I. Allen, "Regularized tensor factorizations and higher-order principal components analysis," arXiv preprint arXiv:1202.2476, 2012.
- [40] G. Allen, "Sparse higher-order principal components analysis.," in AISTATS, vol. 15, 2012.

- [41] Y. Qi, P. Comon, and L.-H. Lim, "Uniqueness of nonnegative tensor approximations," *IEEE Transactions on Information Theory*, vol. 62, no. 4, pp. 2170–2183, 2016.
- [42] C. Lu, J. Feng, Y. Chen, W. Liu, Z. Lin, and S. Yan, "Tensor robust principal component analysis: Exact recovery of corrupted low-rank tensors via convex optimization," in *Proceed*ings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 5249–5257, 2016.
- [43] T. Lesieur, L. Miolane, M. Lelarge, F. Krzakala, and L. Zdeborová, "Statistical and computational phase transitions in spiked tensor estimation," in *Information Theory (ISIT)*, 2017 IEEE International Symposium on, pp. 511–515, IEEE, 2017.
- [44] T. G. Kolda and B. W. Bader, "Tensor decompositions and applications," *SIAM review*, vol. 51, no. 3, pp. 455–500, 2009.
- [45] Q. Berthet and P. Rigollet, "Computational lower bounds for sparse pca," arXiv preprint arXiv:1304.0828, 2013.
- [46] Q. Berthet, P. Rigollet, et al., "Optimal detection of sparse principal components in high dimension," The Annals of Statistics, vol. 41, no. 4, pp. 1780–1815, 2013.
- [47] T. Wang, Q. Berthet, and R. J. Samworth, "Statistical and computational trade-offs in estimation of sparse principal components," arXiv preprint arXiv:1408.5369, 2014.
- [48] C. Gao, Z. Ma, and H. H. Zhou, "Sparse cca: Adaptive estimation and computational barriers," arXiv preprint arXiv:1409.8565, 2014.
- [49] Z. Ma and Y. Wu, "Computational barriers in minimax submatrix detection," *The Annals of Statistics*, vol. 43, no. 3, pp. 1089–1116, 2015.
- [50] T. T. Cai, T. Liang, and A. Rakhlin, "Computational and statistical boundaries for submatrix localization in a large noisy matrix," arXiv preprint arXiv:1502.01988, 2015.
- [51] Y. Chen and J. Xu, "Statistical-computational tradeoffs in planted problems and submatrix localization with a growing number of clusters and submatrices," arXiv preprint arXiv:1402.1267, 2014.

- [52] B. Barak and A. Moitra, "Noisy tensor completion via the sum-of-squares hierarchy," in 29th Annual Conference on Learning Theory, pp. 417–445, 2016.
- [53] B. E. Hajek, Y. Wu, and J. Xu, "Computational lower bounds for community detection on random graphs.," in *COLT*, pp. 899–928, 2015.
- [54] B. Bollobás and P. Erdös, "Cliques in random graphs," in *Mathematical Proceedings of the Cambridge Philosophical Society*, vol. 80, pp. 419–427, Cambridge Univ Press, 1976.
- [55] V. Feldman, E. Grigorescu, L. Reyzin, S. Vempala, and Y. Xiao, "Statistical algorithms and a lower bound for detecting planted cliques," in *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pp. 655–664, ACM, 2013.
- [56] N. Alon, M. Krivelevich, and B. Sudakov, "Finding a large hidden clique in a random graph," *Random Structures and Algorithms*, vol. 13, no. 3-4, pp. 457–466, 1998.
- [57] B. P. Ames and S. A. Vavasis, "Nuclear norm minimization for the planted clique and biclique problems," *Mathematical programming*, vol. 129, no. 1, pp. 69–89, 2011.
- [58] L. Kuvcera, "Expected complexity of graph partitioning problems," *Discrete Applied Mathematics*, vol. 57, no. 2, pp. 193–212, 1995.
- [59] M. Jerrum, "Large cliques elude the metropolis process," Random Structures & Algorithms, vol. 3, no. 4, pp. 347–359, 1992.
- [60] U. Feige and R. Krauthgamer, "The probable value of the lovász–schrijver relaxations for maximum independent set," SIAM Journal on Computing, vol. 32, no. 2, pp. 345–370, 2003.
- [61] P.-A. Wedin, "Perturbation bounds in connection with singular value decomposition," *BIT Numerical Mathematics*, vol. 12, no. 1, pp. 99–111, 1972.
- [62] R. Vershynin, "Introduction to the non-asymptotic analysis of random matrices," arXiv preprint arXiv:1011.3027, 2010.
- [63] V. Koltchinskii and D. Xia, "Optimal estimation of low rank density matrices," Journal of Machine Learning Research, vol. 16, pp. 1757–1792, 2015.

- [64] D. Hush and C. Scovel, "Concentration of the hypergeometric distribution," Statistics & probability letters, vol. 75, no. 2, pp. 127–132, 2005.
- [65] B. Laurent and P. Massart, "Adaptive estimation of a quadratic functional by model selection," *Annals of Statistics*, pp. 1302–1338, 2000.

# Supplement to "Tensor SVD: Statistical and Computational Limits" <sup>1</sup>

Anru Zhang and Dong Xia

University of Wisconsin-Madison

#### Abstract

In this Supplement, we provide additional proofs for the main results and technical lemmas.

# A Additional Proofs

#### A.1 Proof of Theorem 2

We only need to prove upper bounds for  $\hat{U}_k^{\bullet}$  and  $\hat{\mathbf{X}}^{\bullet}$  as the ones for  $\hat{U}_k^{\text{mle}}$  and  $\hat{\mathbf{X}}^{\text{mle}}$  immediately follow. The proof of this theorem is fairly complicated. For convenience, we assume  $\sigma=1$ , denote  $r=\max\{r_1,r_2,r_3\}$ , then  $r\leq C_0p^{1/2}$  according to the assumption. For any orthogonal columns, e.g.  $U_k\in\mathbb{O}_{p_k,r_k}$ , we note  $U_{k\perp}\in\mathbb{O}_{p_k,p_k-r_k}$  as the orthogonal complement of  $U_k$ .

Let  $A \otimes B$  be the Kronecker product between matrices A and B,  $\text{vec}(\cdot)$  be vectorization of matrices and tensors. Similarly as the proof of Theorems 1 and 2 in [29], for any  $V_k \in \mathbb{O}_{p_k, r_k}$ , k = 1, 2, 3,

$$\begin{aligned} & \min_{\hat{\mathbf{S}}} \|\mathbf{Y} - \hat{\mathbf{S}} \times_1 V_1 \times_2 V_2 \times_3 V_3\|_F^2 = \min_{\hat{\mathbf{S}}} \|\text{vec}(\mathbf{Y}) - V_1 \otimes V_2 \otimes V_3 \text{vec}(\hat{\mathbf{S}})\|_2^2 \\ & = \min_{\hat{\mathbf{S}}} \|P_{(V_1 \otimes V_2 \otimes V_3)_{\perp}} \text{vec}(\mathbf{Y}) + P_{V_1 \otimes V_2 \otimes V_3} \text{vec}(\mathbf{Y}) - V_1 \otimes V_2 \otimes V_3 \text{vec}(\hat{\mathbf{S}})\|_2^2 \\ & = \|P_{(V_1 \otimes V_2 \otimes V_3)_{\perp}} \text{vec}(\mathbf{Y})\|_2^2 + \min_{\hat{\mathbf{S}}} \|P_{(V_1 \otimes V_2 \otimes V_3)} \text{vec}(\mathbf{Y}) - V_1 \otimes V_2 \otimes V_3 \text{vec}(\hat{\mathbf{S}})\|_2^2 \\ & = \|P_{(V_1 \otimes V_2 \otimes V_3)_{\perp}} \text{vec}(\mathbf{Y})\|_2^2 = \|\text{vec}(\mathbf{Y})\|_2^2 - \|P_{(V_1 \otimes V_2 \otimes V_3)} \text{vec}(\mathbf{Y})\|_2^2 \\ & = \|\mathbf{Y}\|_F^2 - \|\mathbf{Y} \times_1 V_1^\top \times_2 V_2^\top \times_3 V_3^\top\|_F^2 \end{aligned}$$

<sup>&</sup>lt;sup>1</sup>Anru Zhang is Assistant Professor, Department of Statistics, University of Wisconsin-Madison, Madison, WI 53706, E-mail: anruzhang@stat.wisc.edu; Dong Xia is Visiting Assistant Professor, Department of Statistics, University of Wisconsin-Madison, Madison, WI 53706, E-mail: dongxia@stat.wisc.edu.

where the inequality holds if and only if  $\hat{\mathbf{S}} = \mathbf{Y} \times_1 V_1^{\top} \times_2 V_2^{\top} \times_3 V_3^{\top}$ . Therefore, we must have

$$\|\mathbf{Y} \times_{1} (\hat{U}_{1}^{\bullet})^{\top} \times_{2} (\hat{U}_{2}^{\bullet})^{\top} \times_{3} (\hat{U}_{3}^{\bullet})^{\top}\|_{F}^{2} \geq \|\mathbf{Y} \times_{1} U_{1}^{\top} \times_{2} U_{2}^{\top} \times_{3} U_{3}^{\top}\|_{F}^{2}.$$

For convenience, we simply let  $\hat{U}_k = \hat{U}_k^{\bullet}$  for k = 1, 2, 3 and  $\hat{\mathbf{X}} = \hat{\mathbf{X}}^{\bullet}$  throughout the proof of this theorem. Without loss of generality, we also assume that

$$U_k = \begin{bmatrix} I_{r_k} \\ 0_{(p_k - r_k) \times r_k} \end{bmatrix}, \quad k = 1, 2, 3.$$
 (53)

Such assumption will simplify our notation and make the proof easier to understand. This theorem will be shown by steps.

1. In this first step, we establish some basic probability bounds which will be used in the latter steps. We first let

$$\tilde{\mathbf{X}} = \mathbf{Y}_{[1:r_1,1:r_2,1:r_3]} = \mathbf{S} + \mathbf{Z}_{[1:r_1,1:r_2,1:r_3]} \in \mathbb{R}^{r_1 \times r_2 \times r_3}.$$
 (54)

Then we first have

$$\left\| \tilde{\mathbf{X}} \right\|_{\mathbf{F}} = \left\| \mathbf{Y} \times_1 U_1^{\top} \times_2 U_2^{\top} \times_3 U_3^{\top} \right\|_{\mathbf{F}} \le \left\| \mathbf{Y} \times_1 \hat{U}_1^{\top} \times_2 \hat{U}_2^{\top} \times_3 \hat{U}_3^{\top} \right\|_{\mathbf{F}}. \tag{55}$$

Next, we also note that  $\mathcal{M}_1(\tilde{\mathbf{X}}) = \mathcal{M}_1(\mathbf{S}) + \mathcal{M}_1(\mathbf{Z}_{[1:r_1,1:r_2,1:r_3]}) \in \mathbb{R}^{r_1 \times (r_2 r_3)}$ , i.e., the fixed matrix  $\mathcal{M}_1(\mathbf{S})$  plus an i.i.d. Gaussian matrix. Meanwhile,  $\sigma_{\min}(\mathcal{M}_1(\mathbf{S})) = \sigma_{r_1}(\mathcal{M}_1(\mathbf{X})) \geq \lambda$ . Now, by Lemma 4 in [26],

$$P\left(\sigma_{r_1}^2(\mathcal{M}_1(\tilde{\mathbf{X}})) \ge (\lambda^2 + r_2 r_3)(1 - x)\right) \le C \exp\left(Cr_1 - c(\lambda^2 + r_2 r_3)x^2 \wedge x\right), \quad \text{for any } x > 0.$$

Similar results also hold for  $\sigma_{\min}(\mathcal{M}_2(\tilde{\mathbf{X}}))$  and  $\sigma_{\min}(\mathcal{M}_3(\tilde{\mathbf{X}}))$ . Let x = 1/2, note that  $\lambda \geq C_{gap}p^{1/2}$ , we have

$$\sigma_{\min}(\mathcal{M}_k(\tilde{\mathbf{X}})) \ge \frac{\lambda}{2^{1/2}}, \quad k = 1, 2, 3$$
 (56)

with probability at least  $1 - C \exp(-cp)$  for large enough constant  $C_{gap} > 0$ . Additionally, by Lemma 5, we have

$$\max_{\substack{V_k \in \mathbb{O}_{p_k, r_k} \\ k = 1, 2, 3}} \left\| \mathbf{Z} \times_1 V_1^\top \times_2 V_2^\top \times_3 V_3^\top \right\|_{\mathcal{F}} \le C\sqrt{pr},\tag{57}$$

with probability at least  $1 - C \exp(-cp)$ .

2. In the following Steps 2 and 3, we temporarily ignore the randomness of  $\mathbf{Z}$  and the definition of  $\hat{U}_1, \hat{U}_2, \hat{U}_3$  as the the estimators with high likelihood values. Instead we only assume  $\mathbf{X}, \mathbf{Z} \in \mathbb{R}^{p_1 \times p_2 \times p_3}, \mathbf{X} = \mathbf{S} \times_1 U_1 \times_2 U_2 \times_3 U_3$ , and  $\hat{U}_k \in \mathbb{O}_{p_k, r_k}$  for k = 1, 2, 3 satisfy (55), (56), and (57). By Lemma 1 in [26],  $\|\sin\Theta(\hat{U}_k, U_k)\|_{\mathcal{F}} = \|U_{k\perp}^{\top} \hat{U}_k\|_{\mathcal{F}}$ . Our goal in Steps 2 is to show under such setting, one must has

$$\|\sin\Theta(\hat{U}_k, U_k)\|_{F} = \|U_{k\perp}^{\top} \hat{U}_k\|_{F} \le C \frac{\sqrt{pr}}{\lambda}, \quad k = 1, 2, 3.$$
 (58)

To simplify our notations, we first perform spectral transformation on  $(\hat{U}_k)_{[1:r_k,:]}$  and  $(\hat{U}_k)_{[(r_k+1):p_k,:]}$ . To be specific, let  $(\hat{U}_k)_{[1:r_k,:]} = \bar{U}_k \bar{\Sigma}_k \bar{V}_k^{\top}$  be the singular value decomposition, and  $(\hat{U}_k)_{[(r_k+1):p_k,:]} \bar{V}_k = \bar{Q}_k \bar{R}_k$  be the full QR decomposition (so that  $\bar{Q}_k$  is a square orthogonal matrix). Then we transform

$$\begin{bmatrix} (\hat{U}_k)_{[1:r_k,:]} \\ (\hat{U}_k)_{[(r_k+1):p_k,:]} \end{bmatrix} \Rightarrow \begin{bmatrix} \bar{U}_k^\top \\ \bar{Q}_k^\top \end{bmatrix} \cdot \begin{bmatrix} (\hat{U}_k)_{[1:r_k,:]} \\ (\hat{U}_k)_{[(r_k+1):p,:]} \end{bmatrix} \bar{V}_k = \begin{bmatrix} \bar{U}_k^\top (\hat{U}_k)_{[1:r_k,:]} \bar{V}_k \\ \bar{Q}_k^\top (\hat{U}_k)_{[(r_k+1):p_k,:]} \bar{V}_k \end{bmatrix} = \begin{bmatrix} \bar{\Sigma}_k \\ \bar{R}_k \end{bmatrix},$$

$$\mathbf{Z} \Rightarrow \mathbf{Z} \times_1 \begin{bmatrix} \bar{U}_1 \\ \bar{Q}_1 \end{bmatrix} \times_2 \begin{bmatrix} \bar{U}_2 \\ \bar{Q}_2 \end{bmatrix} \times_3 \begin{bmatrix} \bar{U}_3 \\ \bar{Q}_3 \end{bmatrix},$$

$$\mathbf{X} \Rightarrow \mathbf{X} \times_1 \begin{bmatrix} \bar{U}_1 \\ \bar{Q}_1 \end{bmatrix} \times_2 \begin{bmatrix} \bar{U}_2 \\ \bar{Q}_2 \end{bmatrix} \times_3 \begin{bmatrix} \bar{U}_3 \\ \bar{Q}_3 \end{bmatrix}.$$

We can check that (55), (56), (57) still hold after this transformation. Suppose  $\operatorname{diag}(R_1) = (a_1, \ldots, a_r)$ . Since  $\bar{\Sigma}_1$  is diagonal and  $\bar{R}_1$  is upper diagonal,  $\begin{bmatrix} \bar{\Sigma}_1 \\ \bar{R}_1 \end{bmatrix}$  is orthogonal, we must have all off-diagonal entries of  $\bar{R}_k$  are zero, and  $\bar{\Sigma}_1 = \operatorname{diag}(\sqrt{1 - a_1^2}, \ldots, \sqrt{1 - a_r^2})$ . For convenient we also denote

$$a_i^{(0)} = \sqrt{1 - a_i^2}, a_i^{(1)} = a_i, b_j^{(0)} = \sqrt{1 - b_j^2}, b_j^{(1)} = b_j, c_k^{(0)} = \sqrt{1 - c_k^2}, c_k^{(1)} = c_k.$$
 (59)

Therefore, without loss of generality we can assume there exist real numbers  $0 \le a_i, b_j, c_k \le 1$ 

such that

$$\hat{U}_{1} = \begin{bmatrix} a_{1}^{(0)} & & & & \\ & \ddots & & & \\ & & a_{r_{1}}^{(0)} & & \\ & & \ddots & & \\ & & a_{r_{1}}^{(1)} & & \\ & & \ddots & & \\ & & & a_{r_{1}}^{(1)} \end{bmatrix}, \hat{U}_{2} = \begin{bmatrix} b_{1}^{(0)} & & & & \\ & & \ddots & & \\ & & & b_{r_{2}}^{(0)} & & \\ & & \ddots & & \\ & & & b_{r_{2}}^{(1)} & & \\ & & & \ddots & & \\ & & & b_{r_{2}}^{(1)} & & \\ & & & \ddots & \\ & & & & b_{r_{2}}^{(1)} & \\ & & & & b_{r_{2}}^{(1)} & \\ & & & & b_{r_{2}}^{(1)} & \\ & & & & b_{r_{3}-2r_{3},r_{3}} \end{bmatrix}.$$

$$(60)$$

where  $0_{p_k-2r_k,r_k}$  represents the zero matrix with dimension  $(p_k-2r_k)$ -by- $r_k$ . By the form of  $U_1, U_2, U_3$  in (53), we must have

$$\max_{k=1,2,3} \|\sin\Theta(\hat{U}_k, U_k)\|_{\mathcal{F}} = \max \left\{ (\sum_{i=1}^{r_1} a_i^2)^{1/2}, (\sum_{i=1}^{r_2} b_i^2)^{1/2}, (\sum_{i=1}^{r_3} c_i^2)^{1/2} \right\}.$$
 (61)

In order to show (58) we only need to prove that  $\max\left\{\left(\sum_{i=1}^{r_1}a_i^2\right)^{1/2},\left(\sum_{i=1}^{r_2}b_i^2\right)^{1/2},\left(\sum_{i=1}^{r_3}c_i^2\right)^{1/2}\right\} \leq C\sqrt{pr}$ .

Next, we decompose the noise tensor  $\mathbf{Z}$  to the following eight pieces,

$$\mathbf{Z}^{(t_1 t_2 t_3)} = \mathbf{Z}_{[(t_1 r_1 + 1):(t_1 r_1 + r_1),(t_2 r_2 + 1):(t_2 r_2 + r_2),(t_3 r_3 + 1):(t_3 r_3 + r_3)]}, \quad t_1, t_2, t_3 \in \{0, 1\}.$$
 (62)

By (57), we know  $\|\mathbf{Z}^{(t_1t_2t_3)}\|_{\mathrm{F}} \leq C\sqrt{pr}$ ,  $t_1, t_2, t_3 \in \{0, 1\}$ . Based on the form of  $\hat{U}_1, \hat{U}_2, \hat{U}_3$  (60), we have

$$\begin{split} & \left(\mathbf{Y} \times_{1} \hat{U}_{1} \times_{2} \hat{U}_{2} \times_{3} \hat{U}_{3}\right)_{ijk} \\ = & \tilde{X}_{ijk} a_{i}^{(0)} b_{j}^{(0)} c_{k}^{(0)} + Z_{ijk}^{(100)} a_{i}^{(1)} b_{j}^{(0)} c_{k}^{(0)} + Z_{ijk}^{(010)} a_{i}^{(0)} b_{j}^{(1)} c_{k}^{(0)} + Z_{ijk}^{(001)} a_{i}^{(0)} b_{j}^{(0)} c_{k}^{(1)} \\ & + Z_{ijk}^{(110)} a_{i}^{(1)} b_{j}^{(1)} c_{k}^{(0)} + Z_{ijk}^{(101)} a_{i}^{(1)} b_{j}^{(0)} c_{k}^{(1)} + Z_{ijk}^{(011)} a_{i}^{(0)} b_{j}^{(1)} c_{k}^{(1)} + Z_{ijk}^{(111)} a_{i}^{(1)} b_{j}^{(1)} c_{k}^{(1)} \\ = & \tilde{X}_{ijk} a_{i}^{(0)} b_{j}^{(0)} c_{k}^{(0)} + \sum_{\substack{t_{1}, t_{2}, t_{3} \in \{0, 1\}\\t_{1}, t_{2}, t_{3} \text{ are not all } 0}} Z_{ijk}^{(t_{1}t_{2}t_{3})} a_{i}^{(t_{1})} b_{j}^{(t_{2})} c_{k}^{(t_{3})}. \end{split}$$

Therefore,

$$\begin{split} 0 \overset{(55)}{\leq} & \left\| \mathbf{Y} \times_{1} \hat{U}_{1} \times_{2} \hat{U}_{2} \times_{3} \hat{U}_{3} \right\|_{\mathrm{F}}^{2} - \left\| \mathbf{Y} \times_{1} U_{1} \times_{2} U_{2} \times_{3} U_{3} \right\|_{\mathrm{F}}^{2} \\ &= \sum_{i,j,k=1}^{r_{1},r_{2},r_{3}} \left[ \left( \tilde{X}_{ijk} a_{i}^{(0)} b_{j}^{(0)} c_{k}^{(0)} + \sum_{\substack{t_{1},t_{2},t_{3} \in \{0,1\}\\t_{1},t_{2},t_{3} \text{ are not all } 0}} Z_{ijk}^{(t_{1}t_{2}t_{3})} a_{i}^{(t_{1})} b_{j}^{(t_{2})} c_{k}^{(t_{3})} \right)^{2} - \tilde{X}_{ijk}^{2} \right] \\ &\leq \sum_{i,j,k=1}^{r_{1},r_{2},r_{3}} \tilde{X}_{ijk}^{2} \left[ \left( a_{i}^{(0)} b_{j}^{(0)} c_{k}^{(0)} \right)^{2} - 1 \right] \\ &+ 2 \sum_{\substack{t_{1},t_{2},t_{3} \in \{0,1\}\\t_{1},t_{2},t_{3} \text{ are not all } 0}} \sum_{i,j,k=1}^{r_{1},r_{2},r_{3}} \tilde{X}_{ijk} \left( a_{i}^{(0)} b_{j}^{(0)} c_{k}^{(0)} \right)^{2} Z_{ijk}^{(t_{1}t_{2}t_{3})} a_{i}^{(t_{1})} b_{j}^{(t_{2})} c_{k}^{(t_{3})} \\ &+ 7 \sum_{\substack{t_{1},t_{2},t_{3} \in \{0,1\}\\t_{1},t_{2},t_{3} \text{ are not all } 0}} \sum_{i,j,k=1}^{r_{1},r_{2},r_{3}} \left( Z_{ijk}^{(t_{1}t_{2}t_{3})} a_{i}^{(t_{1})} b_{j}^{(t_{2})} c_{k}^{(t_{3})} \right)^{2} \\ &\leq \sum_{i,j,k=1}^{r_{1},r_{2},r_{3}} \tilde{X}_{ijk}^{2} \left[ (1-a_{i}^{2})(1-b_{j}^{2})(1-c_{k}^{2}) - 1 \right] \\ &+ 63 \sum_{\substack{t_{1},t_{2},t_{3} \in \{0,1\}\\t_{1},t_{2},t_{3} \in \text{ en ot all } 0}} \sum_{i,j,k=1}^{r_{1},r_{2},r_{3}} \left\{ \tilde{X}_{ijk} \sqrt{1-a_{i}^{2}} \sqrt{1-b_{j}^{2}} \sqrt{1-c_{k}^{2}} Z_{ijk}^{(t_{1}t_{2}t_{3})} a_{i}^{(t_{1})} b_{j}^{(t_{2})} c_{k}^{(t_{3})} \right)^{2} \right\}. \end{split}$$

By the inequality above, one of the following inequalities must hold for some  $t_1, t_2, t_3 \in \{0, 1\}$  and  $t_1, t_2, t_3$  are not all 0:

$$0 \leq \sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 \left[ (1 - a_i^2)(1 - b_j^2)(1 - c_k^2) - 1 \right] + 63\tilde{X}_{ijk} \sqrt{1 - a_i^2} \sqrt{1 - b_j^2} \sqrt{1 - c_k^2} Z_{ijk}^{(t_1 t_2 t_3)} a_i^{(t_1)} b_j^{(t_2)} c_k^{(t_3)},$$

$$(63)$$

$$0 \le \sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 \left[ (1-a_i^2)(1-b_j^2)(1-c_k^2) - 1 \right] + 63 \sum_{i,j,k=1}^{r_1,r_2,r_3} \left( Z_{ijk}^{(t_1t_2t_3)} a_i^{(t_1)} b_j^{(t_2)} c_k^{(t_3)} \right)^2.$$
 (64)

- 3. Next we discuss in two different situations: (63) or (64) hold.
  - (a) When (63) holds, we first assume  $t_1 = 1, t_2 = t_3 = 0$  as the other situations follow

similarly. Then

$$0 \leq \sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 \left[ (1-a_i^2) - 1 \right] + 63 \sum_{i,j,k=1}^{r_1,r_2,r_3} |\tilde{X}_{ijk}| |Z_{ijk}^{(100)}| a_i \quad \text{(since } 1 \leq a_i, b_j, c_k \leq 1)$$

$$\leq -\sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 a_i^2 + 63 \left( \sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 a_i^2 \right)^{1/2} \left( \sum_{i,j,k=1}^{r_1,r_2,r_3} (Z_{ijk}^{(100)})^2 \right)^{1/2}, \quad \text{(Cauchy-Schwarz inequality)}$$

Thus,

$$\sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 a_i^2 \le C \sum_{i,j,k=1}^{r_1,r_2,r_3} (Z_{ijk}^{(100)})^2.$$
(65)

Additionally, since

$$\sum_{i,k=1}^{r_2,r_3} \tilde{X}_{ijk}^2 = \left\| \left( \mathcal{M}_1(\tilde{\mathbf{X}}) \right)_{[i,:]} \right\|_2 \ge \sigma_{\min}^2(\mathcal{M}_1(\tilde{\mathbf{X}}))$$

and (56), we have

$$\sum_{i=1}^{r_1} a_i^2 \frac{\lambda^2}{2} \le \sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 a_i^2 \le 63 \sum_{i,j,k=1}^{r_1,r_2,r_3} Z_{ijk}^{(100)} = 63 \|\mathbf{Z}^{(100)}\|_{\mathrm{F}}^2 \le Cpr,$$

which means  $\|\sin\Theta(\hat{U}_1,U_1)\|_{\mathrm{F}} \stackrel{(61)}{=} \sqrt{\sum_{i=1}^{r_1} a_i^2} \leq C\sqrt{pr}/\lambda$ . On the other hand, by (63),

$$0 \le \sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 [(1-a_i^2)(1-b_j^2) - 1] + 63 \sum_{i,j,k=1}^{r_1,r_2,r_3} |\tilde{X}_{ijk}| |Z_{ijk}^{(100)}| a_i \sqrt{1-a_i^2} (1-b_j^2)$$

(by Algorithmic-geometric inequality)

$$\leq \sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 \left[ -b_j^2 - a_i^2 + a_i^2 b_j^2 \right] + \sum_{i,j,k=1}^{r_1,r_2,r_3} \left( \tilde{X}_{ijk}^2 a_i^2 (1 - b_j^2) + \frac{63^2}{4} (Z_{ijk}^{(100)})^2 (1 - a_i^2) (1 - b_j^2) \right) \\
\leq -\sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 b_j^2 + \frac{63^2}{4} \sum_{i,j,k}^{r} (Z_{ijk}^{(100)})^2 \stackrel{(56)}{\leq} -\frac{\lambda^2}{2} \sum_{j=1}^{r_2} b_j^2 + \frac{63^2}{4} \|Z^{(100)}\|_F^2. \tag{66}$$

Therefore,

$$\|\sin\Theta(\hat{U}_2, U_2)\|_{\mathcal{F}} \stackrel{(61)}{=} \sqrt{\sum_{j=1}^{r_2} b_j^2} \le C \|\mathbf{Z}^{(100)}\|_{\mathcal{F}} / \lambda \stackrel{(57)}{\le} C \sqrt{pr} / \lambda.$$

By symmetry, one can also show that  $\|\sin\Theta(\hat{U}_3,U_3)\|_{\mathrm{F}} \leq C\sqrt{pr}/\lambda$ . In summary, we must have

$$\max_{k=1,2,3} \|\sin\Theta(\hat{U}_k, U_k)\|_{\mathcal{F}} \le C \frac{\sqrt{pr}}{\lambda}$$

for some constant C > 0 when (63) holds.

(b) When (64) holds for some  $t_1, t_2, t_3 \in \{0, 1\}$ ,

$$\begin{split} 0 &\leq \sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 [(1-a_i^2)(1-b_j^2)(1-c_k^2)-1] + 63 \sum_{i,j,k=1}^{r_1,r_2,r_3} (Z_{ijk}^{(t_1t_2t_3)})^2 (a_i^{(t_1)})^2 (b_j^{(t_2)})^2 (c_k^{(t_3)})^2 \\ &\leq \sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 [(1-a_i^2)-1] + 63 \sum_{i,j,k=1}^{r_1,r_2,r_3} (Z_{ijk}^{(t_1t_2t_3)})^2 \\ &\leq - \sum_{i,j,k}^{r_1,r_2,r_3} a_i^2 \tilde{X}_{ijk}^2 + 63 \|\mathbf{Z}^{(t_1t_2t_3)}\|_{\mathrm{F}}^2 \leq \frac{\lambda^2}{2} \sum_{i,j,k=1}^{r_1,r_2,r_3} a_i^2 + Cpr, \end{split}$$

which means  $\|\sin\Theta(\hat{U}_1,U_1)\|_{\mathrm{F}} \stackrel{\text{(61)}}{=} \sqrt{\sum_{i=1}^{r_1} a_i^2} \leq C\sqrt{pr}/\lambda$ . One can similarly prove the parallel results for  $\|\sin\Theta(\hat{U}_2,U_2)\|_{\mathrm{F}}$  and  $\|\sin\Theta(\hat{U}_3,U_3)\|_{\mathrm{F}}$ . In summary, we also have

$$\max_{k=1,2,3} \|\sin\Theta(\hat{U}_k, U_k)\|_{\mathcal{F}} \le C \frac{\sqrt{pr}}{\lambda}$$

for some constant C > 0 when (64) holds.

To sum up, we have the derived perturbation bound: under (55), (56), (57), one must have  $\|\sin\Theta(\hat{U}_k,U_k)\|_{\text{F}} \leq C\sqrt{pr}/\lambda$ .

4. Next we consider the recovery loss for  $\hat{\mathbf{X}}$ . Similarly as Steps 2-3, we temporarily ignore the randomness of  $\mathbf{Z}$ , and the definition of  $\hat{U}_1, \hat{U}_2, \hat{U}_3$  as the estimators with high likelihood values in this step. We aim to prove

$$\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}} \le C\sqrt{pr},$$

under the assumptions of (55), (56), and (57). First, without loss of generality we can assume  $\hat{U}_1, \hat{U}_2, \hat{U}_3$  have the simple form (60). Based on the structure of  $\hat{U}_1, \hat{U}_2, \hat{U}_3$ , we know

$$P_{\hat{U}_1} = \begin{bmatrix} (a_1^{(0)})^2 & & & a_1^{(1)}a_1^{(0)} \\ & \ddots & & \ddots & & \ddots & 0_{r_1,p_1-2r_1} \\ & & (a_{r_1}^{(0)})^2 & & & a_{r_1}^{(1)}a_{r_1}^{(0)} \\ & & & (a_1^{(0)})^2 & & & a_{r_1}^{(1)}a_{r_1}^{(0)} \\ & & \ddots & & & \ddots & 0_{r_1,p_1-2r_1} \\ & & & a_r^{(1)}a_{r_1}^{(0)} & & & & (a_{r_1}^{(1)})^2 \\ & & & & 0_{p_1-2r_1,r_1} & & 0_{p_1-2r_1,r_1} & & 0_{p_1-2r_1,p_1-2r_1} \end{bmatrix},$$

while  $P_{\hat{U}_2}$  and  $P_{\hat{U}_3}$  can be written in similar forms. We have the following decomposition for

$$\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}}$$

$$\begin{aligned}
& \left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_{F} = \left\| \mathbf{Y} \times_{1} P_{\hat{U}_{1}} \times_{2} P_{\hat{U}_{2}} \times_{3} P_{\hat{U}_{3}} - \mathbf{S} \times_{1} U_{1} \times_{2} U_{2} \times_{3} U_{3} \right\|_{F} \\
& \leq \left\| \mathbf{Z} \times_{1} P_{\hat{U}_{1}} \times_{2} P_{\hat{U}_{2}} \times_{3} P_{\hat{U}_{3}} \right\|_{F} + \left\| \mathbf{X} \times_{1} P_{\hat{U}_{1}} \times_{2} P_{\hat{U}_{2}} \times_{3} P_{\hat{U}_{3}} - \mathbf{S} \times_{1} U_{1} \times_{2} U_{2} \times_{3} U_{3} \right\|_{F} \\
& \leq C \sqrt{pr} + \left\| \mathbf{S} \times_{1} \left( P_{\hat{U}_{1}} U_{1}^{\top} \right) \times_{2} \left( P_{\hat{U}_{2}} U_{2}^{\top} \right) \times_{3} \left( P_{\hat{U}_{3}} U_{3}^{\top} \right) - \mathbf{S} \times_{1} U_{1} \times_{2} U_{2} \times_{3} U_{3} \right\|_{F} \\
& \leq C \sqrt{pr} + \left\| \tilde{\mathbf{X}} \times_{1} \left( P_{\hat{U}_{1}} U_{1}^{\top} \right) \times_{2} \left( P_{\hat{U}_{2}} U_{2}^{\top} \right) \times_{3} \left( P_{\hat{U}_{3}} U_{3}^{\top} \right) - \tilde{\mathbf{X}} \times_{1} U_{1} \times_{2} U_{2} \times_{3} U_{3} \right\|_{F} \\
& + \left\| \mathbf{Z}_{[1:r,1:r,1:r]} \times_{1} \left( P_{\hat{U}_{1}} U_{1}^{\top} \right) \times_{2} \left( P_{\hat{U}_{2}} U_{2}^{\top} \right) \times_{3} \left( P_{\hat{U}_{3}} U_{3}^{\top} \right) \right\| + \left\| \mathbf{Z}_{[1:r,1:r,1:r]} \times_{1} U_{1} \times_{2} U_{2} \times_{3} U_{3} \right\|_{F} \\
& \leq C \sqrt{pr} + \left\| \tilde{\mathbf{X}} \times_{1} \left( P_{\hat{U}_{1}} U_{1}^{\top} \right) \times_{2} \left( P_{\hat{U}_{2}} U_{2}^{\top} \right) \times_{3} \left( P_{\hat{U}_{3}} U_{3}^{\top} \right) - \tilde{\mathbf{X}} \times_{1} U_{1} \times_{2} U_{2} \times_{3} U_{3} \right\|_{F} . \\
& \leq C \sqrt{pr} + \left\| \tilde{\mathbf{X}} \times_{1} \left( P_{\hat{U}_{1}} U_{1}^{\top} \right) \times_{2} \left( P_{\hat{U}_{2}} U_{2}^{\top} \right) \times_{3} \left( P_{\hat{U}_{3}} U_{3}^{\top} \right) - \tilde{\mathbf{X}} \times_{1} U_{1} \times_{2} U_{2} \times_{3} U_{3} \right\|_{F} . \\
& \leq C \sqrt{pr} + \left\| \tilde{\mathbf{X}} \times_{1} \left( P_{\hat{U}_{1}} U_{1}^{\top} \right) \times_{2} \left( P_{\hat{U}_{2}} U_{2}^{\top} \right) \times_{3} \left( P_{\hat{U}_{3}} U_{3}^{\top} \right) - \tilde{\mathbf{X}} \times_{1} U_{1} \times_{2} U_{2} \times_{3} U_{3} \right\|_{F} . \\
& \leq C \sqrt{pr} + \left\| \tilde{\mathbf{X}} \times_{1} \left( P_{\hat{U}_{1}} U_{1}^{\top} \right) \times_{2} \left( P_{\hat{U}_{2}} U_{2}^{\top} \right) \times_{3} \left( P_{\hat{U}_{3}} U_{3}^{\top} \right) - \tilde{\mathbf{X}} \times_{1} U_{1} \times_{2} U_{2} \times_{3} U_{3} \right\|_{F} . \\
& \leq C \sqrt{pr} + \left\| \tilde{\mathbf{X}} \times_{1} \left( P_{\hat{U}_{1}} U_{1}^{\top} \right) \times_{2} \left( P_{\hat{U}_{2}} U_{2}^{\top} \right) \times_{3} \left( P_{\hat{U}_{3}} U_{3}^{\top} \right) - \tilde{\mathbf{X}} \times_{1} U_{1} \times_{2} U_{2} \times_{3} U_{3} \right\|_{F} . \\
& \leq C \sqrt{pr} + \left\| \tilde{\mathbf{X}} \times_{1} \left( P_{\hat{U}_{1}} U_{1}^{\top} \right) \times_{2} \left( P_{\hat{U}_{2}} U_{2}^{\top} \right) \times_{2} \left( P_{\hat{U}_{3}} U_{3}^{\top} \right) - \tilde{\mathbf{X}} \times_{1} U_{1} \times_{2} U_{2} \times_{3} U_{3} \right\|_{F} . \\
& \leq C \sqrt$$

Based on the form of  $U_1, U_2, U_3, \hat{U}_1, \hat{U}_2, \hat{U}_3$ , we have

$$\begin{split} & \left\| \tilde{\mathbf{X}} \times_{1} \left( P_{\hat{U}_{1}} U_{1}^{\top} \right) \times_{2} \left( P_{\hat{U}_{2}} U_{2}^{\top} \right) \times_{3} \left( P_{\hat{U}_{3}} U_{3}^{\top} \right) - \tilde{\mathbf{X}} \times_{1} U_{1} \times_{2} U_{2} \times_{3} U_{3} \right\|_{\mathrm{F}}^{2} \\ &= \sum_{i,j,k=1}^{r_{1},r_{2},r_{3}} \tilde{X}_{ijk}^{2} \left( \left( a_{i}^{(0)} b_{j}^{(0)} c_{k}^{(0)} \right)^{2} - 1 \right)^{2} + \sum_{\substack{t_{1},t_{2},t_{3} \in \{0,1\}\\t_{1},t_{2},t_{3} \text{ are not all } 0}} \sum_{i,j,k=1}^{r_{1},r_{2},r_{3}} \tilde{X}_{ijk}^{2} \left( a_{i}^{(0)} b_{j}^{(0)} c_{k}^{(0)} \right)^{2} \left( \left( a_{i}^{(0)} \right)^{2} + \left( a_{i}^{(1)} \right)^{2} \right) \left( \left( b_{j}^{(0)} \right)^{2} + \left( b_{j}^{(1)} \right)^{2} \right) \left( \left( c_{k}^{(0)} \right)^{2} + \left( c_{k}^{(1)} \right)^{2} \right) \\ &+ \sum_{i,j,k=1}^{r_{1},r_{2},r_{3}} \tilde{X}_{ijk}^{2} \left( -2 \left( a_{i}^{(0)} b_{j}^{(0)} c_{k}^{(0)} \right)^{2} + 1 \right). \end{split}$$

Recall the actual values of  $a_i^{(0)}, a_i^{(1)}, b_j^{(0)}, b_j^{(1)}, c_k^{(0)}, c_k^{(1)}$  in (59), we further have

$$\left\| \tilde{\mathbf{X}} \times_{1} \left( P_{\hat{U}_{1}} U_{1}^{\top} \right) \times_{2} \left( P_{\hat{U}_{2}} U_{2}^{\top} \right) \times_{3} \left( P_{\hat{U}_{3}} U_{3}^{\top} \right) - \tilde{\mathbf{X}} \times_{1} U_{1} \times_{2} U_{2} \times_{3} U_{3} \right\|_{F}^{2} \\
\sum_{i,j,k=1}^{r_{1},r_{2},r_{3}} \tilde{X}_{ijk}^{2} \left( a_{i}^{(0)} b_{j}^{(0)} c_{k}^{(0)} \right)^{2} + \sum_{i,j,k=1}^{r_{1},r_{2},r_{3}} \tilde{X}_{ijk}^{2} \left( -2 \left( a_{i}^{(0)} b_{j}^{(0)} c_{k}^{(0)} \right)^{2} + 1 \right) \\
= \sum_{i,j,k=1}^{r_{1},r_{2},r_{3}} \tilde{X}_{ijk}^{2} \left( 1 - (1 - a_{i}^{2})(1 - b_{j}^{2})(1 - c_{k}^{2}) \right).$$
(68)

By the analysis in Step 2, we know under (55), (56), (57), at least one of (63) and (64) must hold for some  $(t_1, t_2, t_3) \in \{0, 1\}^3 \setminus \{(0, 0, 0)\}$ . Again, we discuss in two different situations to show no matter which of (63) or (63) happen, we must have

$$\sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 \left[ 1 - (1 - a_i^2)(1 - b_j^2)(1 - c_k^2) \right] \le Cpr.$$
 (69)

(a) When (63) holds, we again assume  $t_1 = 1, t_2 = t_3 = 0$  as the other situations follow similarly. Particularly, we have shown in Step 2 (a), (65) and (66),

$$\sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 a_i^2 \leq C \sum_{i,j,k=1}^{r_1,r_2,r_3} (Z_{ijk}^{(100)})^2, \quad \sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 b_j^2 \leq C \sum_{i,j,k=1}^{r_1,r_2,r_3} (Z_{ijk}^{(100)})^2$$

Clearly,  $\sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 c_k^2 \leq C \sum_{i,j,k=1}^{r_1,r_2,r_3} (Z_{ijk}^{(100)})^2$  can be derived by symmetry. Then,

$$\sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 \left(1 - (1 - a_i^2)(1 - b_j^2)(1 - c_k^2)\right)$$

$$= \sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 \left(a_i^2 + b_j^2 + c_k^2 - a_i^2 b_j^2 - a_i^2 c_k^2 - b_j^2 c_k^2 + a_i^2 b_j^2 c_k^2\right)$$

$$\leq \sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 (a_i^2 + b_j^2 + c_k^2) \quad \text{(since } 0 \leq a_i, b_j, c_k \leq 1\text{)}$$

$$\leq C \sum_{i,j,k=1}^{r_1,r_2,r_3} (Z_{ijk}^{(100)})^2 \leq Cpr.$$

(b) When (64) holds, one has

$$\sum_{i,j,k=1}^{r_1,r_2,r_3} \tilde{X}_{ijk}^2 \left(1 - (1 - a_i^2)(1 - b_j^2)(1 - c_k^2)\right) \le 63 \sum_{ijk=1}^{r_1,r_2,r_3} \left(Z_{ijk}^{(t_1t_2t_3)} a_i^{(t_1)} b_j^{(t_2)} c_k^{(t_3)}\right)^2$$

$$\le C \sum_{i,i,k=1}^{r_1,r_2,r_3} \left(Z_{ijk}^{(t_1t_2t_3)}\right)^2 = C \|\mathbf{Z}^{(t_1t_2t_3)}\|_{\mathcal{F}} \le Cpr.$$

In summary of Cases (a)(b), we must have (69). Combining (68), (67), and (69), we have shown

$$\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathcal{F}} \le C\sqrt{pr}$$
, under (55), (56), (57). (70)

5. We finalize the proof for Theorem 2 in this step. We let  $Q = \{(55), (56), (57) \text{ all hold}\}$ . By Step  $1, P(Q) \ge 1 - C \exp(-cp)$ ; by Steps 2-4, one has  $\|\sin\Theta(\hat{U}_k, U_k)\|_F \le C\sqrt{pr}/\lambda, k = 1, 2, 3$ , and  $\|\hat{\mathbf{X}} - \mathbf{X}\|_F \le C\sqrt{pr}$  under Q. The rest of the proof is essentially the same as the Step 4 in the proof of Theorem 1.

Since  $\hat{\mathbf{X}}$  is a projection of  $\mathbf{Y}$  by definition, so  $\|\hat{\mathbf{X}}\|_{F} \leq \|\mathbf{Y}\|_{F} \leq \|\mathbf{X}\|_{F} + \|\mathbf{Z}\|_{F}$ . Then we have the following upper bound for 4-th moment of recovery error,

$$\mathbb{E}\|\hat{\mathbf{X}} - \mathbf{X}\|_{F}^{4} \le C \left(\mathbb{E}\|\hat{\mathbf{X}}\|_{F}^{4} + \|\mathbf{X}\|_{F}^{4}\right) \le C\|\mathbf{X}\|_{F}^{4} + C\mathbb{E}\|\mathbf{Z}\|_{F}^{4}$$
$$\le C \exp(c_{0}p) + C\mathbb{E}\left(\chi_{p^{3}}^{2}\right)^{2} = C \exp(c_{0}p) + Cp^{6}.$$

The we have the following upper bound for the risk of  $\hat{\mathbf{X}}$ ,

$$\mathbb{E}\|\hat{\mathbf{X}} - \mathbf{X}\|_{F}^{2} = \mathbb{E}\|\hat{\mathbf{X}} - \mathbf{X}\|_{F}^{2} \mathbf{1}_{Q} + \mathbb{E}\|\hat{\mathbf{X}} - \mathbf{X}\|_{F}^{2} \mathbf{1}_{Q^{c}} = Cpr + \sqrt{\mathbb{E}\|\hat{\mathbf{X}} - \mathbf{X}\|_{F}^{4} \mathbb{E} \mathbf{1}_{Q^{c}}}$$

$$\leq Cpr + C\exp\left((c_{0} - c)p\right) + Cp^{6}\exp(-cp).$$

Thus, one can select  $c_0 < c$  to ensure that

$$\mathbb{E}\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}}^2 \le Cpr \le p_1r_1 + p_2r_2 + p_3r_3.$$

Additionally, when  $\sigma_{\min}(\mathcal{M}_k(\mathbf{X})) \geq \lambda$ , we have  $\|\mathbf{X}\|_{\mathrm{F}}^2 = \|\mathcal{M}_k(\mathbf{X})\|_{\mathrm{F}}^2 \geq r_k \lambda^2$  for k = 1, 2, 3, which implies  $\|\mathbf{X}\|_{\mathrm{F}}^2 \geq Cr\lambda^2$ . Thus we also have

$$\mathbf{E} \frac{\|\hat{\mathbf{X}} - \mathbf{X}\|_{\mathrm{F}}^2}{\|\mathbf{X}\|_{\mathrm{F}}^2} \le \frac{p_1 + p_2 + p_3}{\lambda^2}.$$

Now we consider the Frobenius  $\sin \theta$  norm risk for  $\hat{U}_k$ . Since  $\sin \Theta(\hat{U}_k, U_k)$  is a  $r_k$ -by- $r_k$  matrix with spectral norm no more than 1, definition  $\|\sin \Theta(\hat{U}_k, U_k)\|_F^2 \le r_k \le r$ . Therefore, one has

$$\mathbb{E}\|\sin\Theta(\hat{U}_k, U_k)\|_{\mathcal{F}} = \mathbb{E}\|\sin\Theta(\hat{U}_k, U_k)\|_{\mathcal{F}} 1_Q + \mathbb{E}\|\sin\Theta(\hat{U}_k, U_k)\|_{\mathcal{F}} 1_{Q^c}$$

$$= C\frac{\sqrt{pr}}{\lambda} + \sqrt{\mathbb{E}\|\sin\Theta(\hat{U}_k, U_k)\|_{\mathcal{F}}^2 \cdot \mathbb{E}1_{Q^c}} \le C\frac{\sqrt{pr}}{\lambda} + \sqrt{r \cdot C \exp(-cp)}$$

By the definition of  $\lambda$ , we know  $\lambda = \sigma_{r_k}(\mathcal{M}_k(\mathbf{X})) \leq \frac{\|\mathbf{X}\|_F}{\sqrt{r_k}} \leq C \frac{\exp(c_0 p)}{\sqrt{r_k}}$ , so we can select  $c_0 > 0$  small enough to ensure that

$$\frac{\sqrt{pr}}{\lambda} \ge \frac{\sqrt{pr^2}}{C \exp(c_0 p)} \ge c\sqrt{r \cdot C \exp(-cp)}.$$

This means  $\mathbb{E}\|\sin\Theta(\hat{U}_k,U_k)\|_{\mathrm{F}} \leq C\frac{\sqrt{pr}}{\lambda}$ . Finally, for any  $1\leq q\leq 2$ , we have

$$\mathbb{E} r_k^{-1/q} \left\| \sin \Theta(\hat{U}_k, U_k) \right\|_q \leq \mathbb{E} r_k^{-1/2} \left\| \sin \Theta(\hat{U}_k, U_k) \right\|_{\mathcal{F}} \leq C \frac{\sqrt{p}}{\lambda} \leq C \frac{\sqrt{p_k}}{\lambda}.$$

To sum up, we have finished the proof for this theorem.  $\square$ 

#### A.2 Proof of Proposition 1

For convenient, we introduce the following notations:  $m_k = |C \cap D_k|$ ,  $u^{(k)} = (1_C)_{D_k}$ ,  $U^{(k)} = u^{(k)}/\|u^{(k)}\|_2$  is the normalized vector of  $u^{(k)}$ ,  $U^{(k)}_{\perp} \in \mathbb{R}^{\frac{N}{6} \times \left(\frac{N^2}{36} - 1\right)}$  is the orthogonal complement

of  $U^{(k)}$ ,  $\tilde{A}_k = \mathcal{M}_k(2 \cdot \mathbf{A}_{[D_1,D_2,D_3]} - 1_{|D_1| \times |D_2| \times |D_3|})$ , for k = 1, 2, 3. Without loss of generality and for convenience of the presentation, we assume N is a multiple of 6. Based on the statement,

$$(\tilde{A}_{1})_{i,p_{3}(j-1)+k} = \begin{cases} 1, & \text{w.p } 1, & \text{if } \left(u_{i}^{(1)}, u_{j}^{(2)}, u_{k}^{(3)}\right) = (1, 1, 1); \\ 1, & \text{w.p } 1/2, & \text{if } \left(u_{i}^{(1)}, u_{j}^{(2)}, u_{k}^{(3)}\right) \neq (1, 1, 1); \\ -1, & \text{w.p } 1.2, & \text{if } \left(u_{i}^{(1)}, u_{j}^{(2)}, u_{k}^{(3)}\right) \neq (1, 1, 1). \end{cases}$$
 (71)

 $A_2$  and  $A_3$  have the similar form. Therefore,  $\tilde{A}_1$  are all 1 in the block of  $(D_1 \cap C) \times ((D_2 \cap C) \otimes (D_3 \cap C))$ , and are with i.i.d. Rademacher entries outside the block. Since C is uniformly randomly selected from  $V_1$ ,  $|V_1| = N/6$ ,  $|D_k| = N/6$ ,  $|C| = \kappa_N$ , we know  $m_1 = |D_1 \cap C|$ ,  $m_2 = |D_1 \cap C|$ ,  $m_3 = |D_1 \cap C|$  satisfy hypergeometric distribution with parameter  $(\kappa_N, N/2, N/6)$ . Based on the concentration inequality of hypergeometric distribution (Theorem 1 in [64]),

$$\frac{\kappa_N}{4} \le m_k = |D_k \cap C| \le \frac{\kappa_N}{2}, \quad k = 1, 2, 3 \tag{72}$$

with probability at least  $1 - C \exp(-c\kappa_N)$ . Now the rest of the proof is similar to Theorem 3 in [26]. By (71), we have

$$\left( (U^{(1)})^{\top} \tilde{A}_1 \right) \in \mathbb{R}^{N^2/36}, \quad \left( (U^{(1)})^{\top} A_1 \right)_j \begin{cases} = \sqrt{m_1}, & j \in ((D_2 \cap C) \otimes (D_3 \cap C)); \\ \sim \frac{W}{\sqrt{m_1}}, & \text{otherwise,} \end{cases}$$

where W has the same distribution as the sum of  $m_1$  i.i.d. Rademacher random variables. Conditioning on C satisfying (72), similarly as the derivation for Equation (1.15) in the Appendix of [26], we can derive

$$\sigma_1^2 \left( (U^{(1)})^\top A_1 \right) \ge \frac{N^2}{36} + \frac{m_1 m_2 m_3}{4}$$
 with probability at least  $1 - C \exp(-cN)$ ; (73)

$$\sigma_2^2(A_1) \le \frac{N^2}{36} + \frac{m_1 m_2 m_3}{8} \quad \text{with probability at least } 1 - C \exp(-cN); \tag{74}$$

$$\|(u_{\perp}^{(1)})^{\top} A_1 P_{(U^{(1)})^{\top} A_1}\| \le C\sqrt{N}$$
 with probability at least  $1 - C \exp(-cN)$ . (75)

Under the circumstance that (72), (73), (73), and (73) all hold, by Proposition 1 in [26], we have

$$\begin{aligned} \|\sin\Theta(\hat{u}_{1}^{\top},U^{(1)})\| &\leq \frac{\sigma_{2}(A_{1})\|(u_{\perp}^{(1)})^{\top}A_{1}P_{(U^{(1)})^{\top}A_{1}}\|}{\sigma_{1}^{2}\left((U^{(1)})^{\top}A_{1}\right) - \sigma_{2}^{2}(A_{1})} \\ &\stackrel{(73)(73)(73)}{\leq} C\frac{\sqrt{N(N^{2} + m_{1}m_{2}m_{3})}}{m_{1}m_{2}m_{3}} \stackrel{(72)}{\leq} C\frac{N^{3/2} + N^{1/2}\kappa_{N}^{3/2}}{\kappa_{N}^{3}}. \end{aligned}$$

Note that  $\liminf_{N\to\infty} \kappa_N/\sqrt{N} = \infty$ ,  $\lim_{N\to\infty} P((72), (73), (73), (73), and (73) all hold) = 1$ , we have

$$\|\sin\Theta(\hat{u}_1^{\top}, (1_C)_{D_1})\| = \|\sin\Theta(\hat{u}_1^{\top}, U^{(1)})\| \stackrel{d}{\to} 0, \text{ as } N \to \infty.$$

The proofs for k=2,3 essentially follow. Therefore, we have finished the proof of this proposition.  $\Box$ 

# B Appendix: Technical Lemmas

We collect all technical lemmas that has been used in the theoretical proofs throughout the paper in this section.

The following lemma shows the equivalence between two widely considered Schatten q-norm distances for singular subspaces.

**Lemma 3.** For any  $U_1, U_2 \in \mathbb{O}_{p,r}$  and all  $1 \leq q \leq +\infty$ ,

$$\frac{1}{4} \|U_1 U_1^{\top} - U_2 U_2^{\top}\|_q \le \|\sin\Theta(U_1, U_2)\|_q \le \|U_1 U_1^{\top} - U_2 U_2^{\top}\|_q.$$

We will use the following properties of tensor algebra in the technical analysis of this paper.

Lemma 4 (Properties in Tensor Algebra).

• Suppose  $\mathbf{X} \in \mathbb{R}^{p_1 \times p_2 \times p_3}$ ,  $U_k \in \mathbb{R}^{p_k \times r_k}$  for k = 1, 2, 3. Then we have the following identity related to tensor matricizations,

$$\mathcal{M}_k \left( \mathbf{X} \times_{k+1} U_{k+1}^{\top} \times_{k+2} U_{k+2}^{\top} \right) = \mathcal{M}_k(\mathbf{X}) (U_{k+1} \otimes U_{k+2}), \quad k = 1, 2, 3.$$
 (76)

• Suppose we further have  $\tilde{U}_k \in \mathbb{R}^{p_k \times \tilde{r}_k}$  for k = 1, 2, 3, then

$$\left(\tilde{U}_{1} \otimes \tilde{U}_{2}\right)^{\top} = \left(\tilde{U}_{1}^{\top}\right) \otimes \left(\tilde{U}_{2}^{\top}\right), \quad \left(\tilde{U}_{2} \otimes \tilde{U}_{3}\right)^{\top} \left(U_{2} \otimes U_{3}\right) = \left(\tilde{U}_{2}^{\top} U_{2}\right) \otimes \left(\tilde{U}_{3}^{\top} U_{3}\right). \tag{77}$$

$$||U_2 \otimes U_3|| = ||U_2|| \cdot ||U_3||, \quad ||U_2 \otimes U_3||_{\mathcal{F}} = ||U_2||_{\mathcal{F}} \cdot ||U_3||_{\mathcal{F}},$$

$$\sigma_{\min}(U_2 \otimes U_3) = \sigma_{\min}(U_2)\sigma_{\min}(U_3).$$
(78)

• (Properties related to projections) Suppose  $U_2 \in \mathbb{O}_{p_2,r_2}, U_3 \in \mathbb{O}_{p_3,r_3}$ , and  $U_{2\perp} \in \mathbb{O}_{p_2,p_2-r_2}, U_{3\perp} \in \mathbb{O}_{p_3,p_3-r_3}$  are their orthogonal complement, respectively. Then  $P_{U_2 \otimes U_3} = P_{U_2} \otimes P_{U_3}$ , and we have the following decomposition

$$I_{p_2p_3} = P_{I_{p_2} \otimes U_3} + P_{I_{p_3} \otimes U_{3\perp}} = P_{U_2 \otimes I_{p_3}} + P_{U_{3\perp} \otimes I_{p_2}}$$

$$= P_{U_2 \otimes U_3} + P_{U_{2\perp} \otimes U_3} + P_{U_2 \otimes U_{3\perp}} + P_{U_{2\perp} \otimes U_{3\perp}}.$$
(79)

The following lemma characterizes the maximum of norms for i.i.d. Gaussian tensors after any projections.

**Lemma 5.** For i.i.d. Gaussian tensor  $\mathbf{Z} \in \mathbb{R}^{p_1 \times p_2 \times p_3}$ ,  $\mathbf{Z} \stackrel{iid}{\sim} N(0,1)$ , we have the following tail bound for the projections,

$$P\left(\max_{\substack{V_2 \in \mathbb{R}^{p_2 \times r_2}, V_3 \in \mathbb{R}^{p_3 \times r_3} \\ \|V_2\| \leq 1, \|V_3\| \leq 1}} \left\| \mathcal{M}_1\left(\mathbf{Z} \times_2 V_2^\top \times_3 V_3^\top\right) \right\| \geq C\sqrt{p_1} + C\sqrt{r_2r_3} + C\sqrt{1+t}\left(\sqrt{p_2r_2} + \sqrt{p_3r_3}\right) \right)$$

$$\leq C \exp(-Ct(p_2r_2 + p_3r_3))$$

for any t > 0. Similar results also hold for  $\mathcal{M}_2\left(\mathbf{Z} \times_1 V_1^{\top} \times_3 V_3^{\top}\right)$  and  $\mathcal{M}_3\left(\mathbf{Z} \times_1 V_1^{\top} \times_2 V_2^{\top}\right)$ . Meanwhile, there exists uniform C > 0 such that

$$P\left(\max_{\substack{V_1,V_2,V_3\in\mathbb{R}^{p\times r}\\\max\{\|V_1\|,\|V_2\|,\|V_3\|\}\leq 1}} \left\|\mathbf{Z}\times_1V_1^{\top}\times_2V_2^{\top}\times_3V_3^{\top}\right\|_{\mathrm{F}}^2 \geq Cr_1r_2r_3 + C(1+t)(p_1r_1 + p_2r_2 + p_3r_3)\right)$$

$$\leq \exp\left(-Ct(p_1r_1 + p_2r_2 + p_3r_3)\right)$$
(80)

for any t > 0.

In the perturbation bound analysis in this paper, we also need the following technical result to bound the spectral and Frobenius norm for the projections.

**Lemma 6.** Suppose  $X, Z \in \mathbb{R}^{p_1 \times p_2}$ , rank(X) = r. If the singular value decomposition of X and Y are written as

$$Y = X + Z = \hat{U}\hat{\Sigma}\hat{V}^{\top} = \begin{bmatrix} \hat{U}_1 & \hat{U}_2 \end{bmatrix} \cdot \begin{bmatrix} \hat{\Sigma}_1 \\ & \hat{\Sigma}_2 \end{bmatrix} \cdot \begin{bmatrix} \hat{V}_1^{\top} & \hat{V}_2^{\top} \end{bmatrix},$$

where  $\hat{U}_1 \in \mathbb{O}_{p_1,r}, \hat{V}_1 \in \mathbb{O}_{p_2,r}$  correspond to the leading r left and right singular vectors; and  $\hat{U}_2 \in \mathbb{O}_{p_1,p_2-r}, \hat{V}_1 \in \mathbb{O}_{p_2,p_2-r}$  correspond to their orthonormal complement. Then

$$\left\|P_{\hat{U}_2}X\right\| \leq 2\|Z\|, \quad \left\|P_{\hat{U}_2}X\right\|_{\mathcal{F}} \leq \min\left\{2\sqrt{r}\|Z\|, 2\|Z\|_{\mathcal{F}}\right\}.$$

The following lemma provides a detailed analysis for  $\varepsilon$ -net for the class of regular matrices under various norms and for the low-rank matrices under spectral norm.

**Lemma 7** ( $\varepsilon$ -net for Regular and Low-rank Matrices).

- Suppose  $\|\cdot\|_{\bullet}$  is any matrix norm,  $\mathcal{X}_{p_1,p_2} = \{X \in \mathbb{R}^{p_1 \times p_2} : \|X\| \leq 1\}$  is the unit ball around the center in  $\|\cdot\|_{\bullet}$  norm. Then there exists an  $\varepsilon$ -net  $\bar{\mathcal{X}}_{p_1,p_2}$  in  $\|\cdot\|_{\bullet}$  norm with cardinality at most  $((2+\varepsilon)/\varepsilon)^{p_1p_2}$  for  $\mathcal{X}_{p_1,p_2}$ . To be specific, there exists  $X^{(1)}, \ldots, X^{(N)}$  with  $N \leq ((2+\varepsilon)/\varepsilon)^{p_1p_2}$ , such that for all  $X \in \mathcal{X}_{p_1,p_2}$ , there exists  $i \in \{1,\ldots,N\}$  satisfying  $\|X^{(i)} X\| \leq \varepsilon$ .
- Let  $\mathcal{X}_{p_1,p_2,r} = \{X \in \mathbb{R}^{p_1 \times p_2} : \operatorname{rank}(X) \leq r, \|X\| \leq 1\}$  be the class of low-rank matrices under spectral norm. Then there exists an  $\varepsilon$ -net  $\bar{\mathcal{X}}_r$  for  $\mathcal{X}_{p_1,p_2,r}$  with cardinality at most  $((4+\varepsilon)/\varepsilon)^{(p_1+p_2)r}$ . Specifically, there exists  $X^{(1)}, \ldots, X^{(N)}$  with  $N \leq ((4+\varepsilon)/\varepsilon)^{(p_1+p_2)r}$ , such that for all  $X \in \mathcal{X}_{p_1,p_2,r}$ , there exists  $i \in \{1,\ldots,N\}$  satisfying  $\|X^{(i)} X\| \leq \varepsilon$ .

The next lemma characterizes the tail probability for i.i.d. Gaussian vector after multiplication of any fixed matrix.

**Lemma 8.** Suppose  $u \in \mathbb{R}^p$  such that  $x \stackrel{iid}{\sim} N(0,1)$ ,  $A \in \mathbb{R}^{p \times r}$  is a fixed matrix, then

$$P\left(\|Au\|_{2}^{2} - \|A\|_{F}^{2} \le -2\|A^{\top}A\|_{F}\sqrt{t}\right) \le \exp(-x);$$

$$P\left(\|Au\|_{2}^{2} - \|A\|_{F}^{2} \ge 2\|A^{\top}A\|_{F}\sqrt{t} + 2\|A\|^{2}t\right) \le \exp(-x).$$

# C Proof of Technical Lemmas

#### C.1 Proof of Lemma 1

Without loss of generality, assume that  $p \equiv 0 \pmod{2}$ . Hereafter, set N = 3p and  $\kappa_N = 20k$  with  $k = \lfloor p^{(1-\tau)/2} \rfloor$ . Our main technique is based on a reduction scheme which maps any adjacency

tensor  $\mathbf{A} \in \{0,1\}^{N \times N \times N}$  to a random tensor  $\mathbf{Y} \in \mathbb{R}^{p \times p \times p}$  in  $O(N^3)$  number of flops. The technique was invented in [49], adapted from a bottom-left trick in [45]. Some other related methods can be found in [50] and [47]. For the completeness and readability of our paper, we provide a detailed application of this technique to the tensor settings.

To this end, for any  $M \geq 3$  and  $0 < \mu \leq \frac{1}{2M}$ , define two random variables

$$\xi^+ := (Z + \mu) \mathbf{1}(|Z| \le M)$$
 and  $\xi^- := (\tilde{Z} - \mu) \mathbf{1}(|\tilde{Z}| \le M)$ 

where Z and  $\tilde{Z}$  denote independent standard normal random variables. The randomized mapping from  $\mathbf{A} \in \{0,1\}^{N \times N \times N}$  to a random matrix  $\mathbf{Y} \in \mathbb{R}^{p \times p \times p}$  is essentially one step of Gaussianization. For simplicity, denote  $V_1 := \{1,2,\ldots,\frac{p}{2}\} \cup \{\frac{3p}{2}+1,\ldots,2p\}$ ,

$$V_2 := \left\{ \frac{p}{2} + 1, \dots, p \right\} \cup \left\{ 2p + 1, \dots, \frac{5p}{2} \right\},$$

and

$$V_3 := \{p+1, \dots, \frac{3p}{2}\} \cup \{\frac{5p}{2} + 1, \dots, 3p\}.$$

Therefore,  $V_1, V_2, V_3$  are disjoint and  $V_1 \cup V_2 \cup V_3 = [N]$ . Given an adjacency tensor  $\mathbf{A} \in \{0,1\}^{N \times N \times N}$ , let  $\mathbf{A}_0 = \mathbf{A}_{V_1,V_2,V_3} \in \mathbb{R}^{p \times p \times p}$  be a corner block of  $\mathbf{A}$ . Conditioned on  $\mathbf{A}_0$ , we generate a random tensor  $\mathbf{Y} \in \mathbb{R}^{p \times p \times p}$  such that

$$Y_{a,b,c} = (1 - (A_0)_{a,b,c}) \Xi_{a,b,c}^- + (A_0)_{a,b,c} \Xi_{a,b,c}^+, \quad \forall a, b, c \in [p]$$

where  $\mathbf{\Xi}^- \in \mathbb{R}^{p \times p \times p}$  has i.i.d. entries with the same distribution as  $\xi^-$  and  $\mathbf{\Xi}^+ \in \mathbb{R}^{p \times p \times p}$  has i.i.d. entries with the same distribution as  $\xi^+$ . Clearly, this process defines a deterministic map for any fixed  $\mathbf{\Xi}^-, \mathbf{\Xi}^+ \in \mathbb{R}^{p \times p \times p}$ 

$$\mathcal{T}: \{0,1\}^{N \times N \times N} \times \mathbb{R}^{p \times p \times p} \times \mathbb{R}^{p \times p \times p} \mapsto \mathbb{R}^{p \times p \times p}$$
$$(\mathbf{A}, \mathbf{\Xi}^-, \mathbf{\Xi}^+) \mapsto \mathbf{Y}.$$

Let  $\mathcal{L}(\mathbf{X})$  denote the law of a random tensor  $\mathbf{X}$ . The total variation distance between two probability distributions  $\mathbb{P}_1$  and  $\mathbb{P}_2$  is denoted by  $d_{\mathrm{TV}}(\mathbb{P}_1, \mathbb{P}_2)$ . The following lemma is analogous to [49, Lemma 2] and the proof is skipped here.

**Lemma 9.** Let  $M \ge 4$ ,  $\mu \le \frac{1}{2M}$ , and  $\eta$  be a Bernoulli random variable. Suppose  $\xi$  is a random variable such that  $(\xi | \eta = 1) = \xi^+$  and  $(\xi | \eta = 0) = \xi^-$ .

(1) If 
$$\mathbb{P}(\eta = 1) = 1$$
, then  $d_{\text{TV}}(\mathcal{L}(\xi), \mathcal{N}(\mu, 1)) \le e^{(1-M^2)/2}$ ;

(2) If 
$$\mathbb{P}(\eta = 0) = \mathbb{P}(\eta = 1) = \frac{1}{2}$$
, then  $d_{\text{TV}}(\mathcal{L}(\xi), \mathcal{N}(0, 1)) \le e^{-M^2/2}$ .

Our next step is to show that (by choosing  $M = \sqrt{8 \log 3p}$  and  $\mu = (2M)^{-1}$ ), the law of  $\mathbf{Y} = \mathcal{T}(\mathbf{A}, \mathbf{\Xi}^-, \mathbf{\Xi}^+)$  is asymptotically equivalent to a mixture over  $\{\mathbb{P}_{\mathbf{X}} : \mathbf{X} \in \mathcal{M}_0(\boldsymbol{p}, k, \boldsymbol{r}, \lambda)\}$  for  $\lambda = \frac{p^{3(1-\tau)/4}}{2\sqrt{8 \log 3p}}$  if  $G \sim H_0$ . On the other hand, if  $G \sim H_1$ , the law of  $\mathbf{Y} = \mathcal{T}(\mathbf{A}, \mathbf{\Xi}^-, \mathbf{\Xi}^+)$  is asymptotically equivalent to a mixture over  $\{\mathbb{P}_{\mathbf{X}} : \mathbf{X} \in \mathcal{M}_1(\boldsymbol{p}, k, \boldsymbol{r}, \lambda)\}$ . For an adjacency tensor  $\mathbf{A} \in \{0, 1\}^{N \times N \times N}$ , we have  $\mathcal{T}(\mathbf{A}, \mathbf{\Xi}^-, \mathbf{\Xi}^+) \in \mathbb{R}^{p \times p \times p}$ . Recall that  $\mathbf{Y} \in \mathbb{R}^{p_1 \times p_2 \times p_3}$  and we define an embedding  $\ell : \mathbb{R}^{p \times p \times p} \mapsto \mathbb{R}^{p_1 \times p_2 \times p_3}$ ,

$$\ell(\mathbf{A})_{ijk} = \begin{cases} A_{ijk} & \text{if } (i,j,k) \in [p] \times [p] \times [p]; \\ 0 & \text{otherwise.} \end{cases}$$
(81)

Lemma 10 is similar to [49, Lemma 4]. We postpone the proof of Lemma 10 to the Appendix.

**Lemma 10.** Let  $\mathbf{A} \in \mathbb{R}^{N \times N \times N}$  be the adjacency tensor of a hypergraph G sampled from either  $H_0$  or  $H_1$  and  $\mathbf{Y} = \ell \circ \mathcal{T}(\mathbf{A}, \mathbf{\Xi}^-, \mathbf{\Xi}^+)$ . Suppose that  $M = \sqrt{8 \log N}$  and  $\mu = \frac{1}{2M}$ . For each i = 0, 1, if  $G \sim H_i$ , there exists a prior distribution  $\pi_i$  on  $\mathcal{M}_i(\mathbf{p}, k, \mathbf{r}, \lambda)$  with  $\lambda = \frac{p^{3/4(1-\tau)}}{2\sqrt{8 \log 3p}}$  such that

$$d_{\text{TV}}(\mathcal{L}(\mathbf{Y}), \mathbb{P}_{\pi_i}) \le \frac{\sqrt{e}}{27N} + 6k(0.86)^{2.5k}.$$

where  $\mathbb{P}_{\pi_i} = \int_{\mathcal{M}_i(\boldsymbol{p},k,\boldsymbol{r},\lambda)} \mathbb{P}_{\mathbf{X}}(\cdot)\pi_i(d\mathbf{X}).$ 

Now, on the contradictory, suppose that the claim of Lemma 1 does not hold. It means that there exists a sequence of polynomial-time tests  $\{\phi_{p_t}\}$  with a sub-sequence  $(p_t)_{t=1}^{\infty}$  of positive integers such that

$$\lim_{t \to \infty} \mathcal{R}_{\boldsymbol{p}, \boldsymbol{r}, \lambda}(\phi_{p_t}) = \lim_{t \to \infty} \left\{ \sup_{\mathbf{X} \in \mathcal{M}_0(\boldsymbol{p}, k, \boldsymbol{r}, \lambda)} \mathbb{P}_{\mathbf{X}} \left\{ \phi_{p_t}(\mathbf{Y}) = 1 \right\} + \sup_{\mathbf{X} \in \mathcal{M}_1(\boldsymbol{p}, k, \boldsymbol{r}, \lambda)} \mathbb{P}_{\mathbf{X}} \left\{ \phi_{p_t}(\mathbf{Y}) = 0 \right\} \right\} < \frac{1}{2}.$$

Define the test  $\psi_{N_t}(\mathbf{A}) = \phi_{p_t}(\ell \circ \mathcal{T}(\mathbf{A}, \mathbf{\Xi}^-, \mathbf{\Xi}^+))$  and we obtain a sequence of polynomial-time tests  $\{\psi_{N_t}\}$  for problem (16) with  $N_t = 3p_t$  for  $t = 1, \dots, \infty$ . It suffices to compute

$$\mathcal{R}_{N_t,\kappa_{N_t}}(\psi_{N_t}) = \mathbb{P}_{H_0}\{\psi_{N_t}(\mathbf{A}) = 1\} + \mathbb{P}_{H_1}\{\psi_{N_t}(\mathbf{A}) = 0\}$$

with  $\kappa_{N_t} = 20 \lfloor p_t^{(1-\tau)/2} \rfloor$ . Note that  $\lim_{t\to\infty} \frac{\log \kappa_{N_t}}{\log \sqrt{N_t}} \le 1-\tau$ . By definition of  $d_{\text{TV}}$  and Lemma 10, under  $H_0$ ,

$$\left| \mathbb{P}_{H_0} \left\{ \psi_{N_t}(\mathbf{A}) = 1 \right\} - \mathbb{P}_{\pi_0} \left\{ \phi_{p_t}(\mathbf{Y}) = 1 \right\} \right|$$

$$= \left| \mathbb{P}_{H_0} \left\{ \phi_{p_t} \left( \ell \circ \mathcal{T}(\mathbf{A}, \mathbf{\Xi}^-, \mathbf{\Xi}^+) \right) = 1 \right\} - \mathbb{P}_{\pi_0} \left\{ \phi_{p_t}(\mathbf{Y}) = 1 \right\} \right|$$

$$\leq d_{\text{TV}} \left( \mathcal{L} \left( \ell \circ \mathcal{T}(\mathbf{A}, \mathbf{\Xi}^-, \mathbf{\Xi}^+) \right), \mathbb{P}_{\pi_0} \right) \leq \frac{\sqrt{e}}{27N_t} + 6k_t (0.86)^{2.5k_t}$$

where  $k_t = \lfloor p_t^{(1-\tau)/2} \rfloor$  and we used the fact that the mixture  $\mathbb{P}_{\pi_0}$  over  $\mathcal{M}_0(\boldsymbol{p}, k, \boldsymbol{r}, \lambda)$  is also a mixture over  $\mathcal{F}_{\boldsymbol{p}, \boldsymbol{r}}(\lambda)$ . In a similar fashion,

$$\left| \mathbb{P}_{H_1} \left\{ \psi_{N_t}(\mathbf{A}) = 0 \right\} - \mathbb{P}_{\pi_1} \left\{ \phi_{p_t}(\mathbf{Y}) = 0 \right\} \right| \le \frac{\sqrt{e}}{27N_t} + 6k_t(0.86)^{2.5k_t}$$

As a result,

$$\mathcal{R}_{N_t,\kappa_{N_t}}(\psi_{N_t}) = \mathbb{P}_{H_0}\left\{\psi_{N_t}(\mathbf{A}) = 1\right\} + \mathbb{P}_{H_1}\left\{\psi_{N_t}(\mathbf{A}) = 0\right\}$$

$$\leq \mathbb{P}_{\pi_0}\left\{\phi_{p_t}(\mathbf{Y}) = 1\right\} + \mathbb{P}_{\pi_1}\left\{\phi_{p_t}(\mathbf{Y}) = 0\right\} + \frac{2\sqrt{e}}{27N_t} + 12k_t(0.86)^{2.5k_t}$$

$$\leq \sup_{\mathbf{X}\in\mathcal{M}_0(\mathbf{p},k,\mathbf{r},\lambda)} \mathbb{P}_{\mathbf{X}}\left\{\phi_{p_t}(\mathbf{Y}) = 1\right\} + \sup_{\mathbf{X}\in\mathcal{M}_1(\mathbf{p},k,\mathbf{r},\lambda)} \mathbb{P}_{\mathbf{X}}\left\{\phi_{p_t}(\mathbf{Y}) = 0\right\} + \frac{2\sqrt{e}}{27N_t} + 12k_t(0.86)^{2.5k_t}$$

Therefore,

$$\lim_{t \to \infty} \mathcal{R}_{N_t, \kappa_{N_t}}(\psi_{N_t}) < \frac{1}{2},$$

which contradicts the hypothesis  $\mathbf{H}(\tau)$ .

## C.2 Proof of Lemma 3

Let  $\sigma_1 \geq \sigma_2 \geq \ldots \geq \sigma_r$  denote the singular values of  $U_1^\top U_2$ . It is easy to check that the singular values of  $U_{1\perp}^\top U_2$  are  $\sqrt{1-\sigma_r^2} \geq \ldots \geq \sqrt{1-\sigma_2^2} \geq \sqrt{1-\sigma_1^2}$ , in view of the fact

$$(U_1^{\top}U_2)^{\top}(U_1^{\top}U_2) + (U_{1\perp}^{\top}U_2)^{\top}(U_{1\perp}^{\top}U_2) = U_2^{\top}U_2 = I_r.$$

Recall that for all  $1 \le q \le +\infty$ ,

$$\|\sin\Theta(U_1, U_2)\|_q = \left(\sum_{i=1}^r \left(\sin(\cos^{-1}\sigma_i)\right)^q\right)^{1/q} = \left(\sum_{i=1}^r \left(1 - \sigma_i^2\right)^{q/2}\right)^{1/q}.$$

The following fact is straightforward:

$$\|U_2 U_2^{\top} - U_1 U_1^{\top}\|_q \ge \|U_{1\perp}^{\top} U_2 U_2^{\top}\|_q = \|U_{1\perp}^{\top} U_2\|_q = \left(\sum_{i=1}^r \left(1 - \sigma_i^2\right)^{q/2}\right)^{1/q}$$

which concludes  $||U_2U_2^{\top} - U_1U_1^{\top}||_q \ge ||\sin\Theta(U_1, U_2)||_q$ . On the other hand,

$$||U_2U_2^{\top} - U_1U_1^{\top}||_q \le ||P_{U_1}(U_2U_2^{\top} - U_1U_1^{\top})P_{U_1}||_q + ||P_{U_1}(U_2U_2^{\top})P_{U_1}^{\perp}||_q$$

$$+ \|P_{U_{1}}^{\perp}(U_{2}U_{2}^{\top})P_{U_{1}}\|_{q} + \|P_{U_{1}}^{\perp}(U_{2}U_{2}^{\top})P_{U_{1}}^{\perp}\|_{q}$$

$$\leq \|U_{1}(U_{1}^{\top}U_{2}U_{2}^{\top}U_{1} - I_{r})U_{1}^{\top}\|_{q} + \|U_{2}^{\top}U_{1\perp}\|_{q} + \|U_{1\perp}^{\top}U_{2}\|_{q} + \|U_{1\perp}^{\top}U_{2}U_{2}^{\top}U_{1\perp}\|_{q}$$

$$\leq 4\left(\sum_{i=1}^{r} \left(1 - \sigma_{i}^{2}\right)^{q/2}\right)^{1/q} \leq 4\|\sin\Theta(U_{1}, U_{2})\|_{q}$$

where we used the fact  $1 - \sigma_i^2 \le \sqrt{1 - \sigma_i^2}$  for all  $1 \le i \le r$ .  $\square$ 

#### C.3 Proof of Lemma 4

• First, we shall note that both  $\mathcal{M}_k\left(\mathbf{X}\times_{k+1}U_{k+1}^{\top}\times_{k+2}U_{k+2}^{\top}\right)$  and  $\mathcal{M}_k(\mathbf{X})(U_{k+1}\otimes U_{k+2})$  are of dimension  $p_k$ -by- $(r_{k+1}r_{k+2})$ . To prove they are equal, we just need to compare each of their entries. We focus on k=1 as the k=2,3 essentially follows. For any  $1\leq i_1\leq p_1, 1\leq i_2\leq r_2, 1\leq i_3\leq r_3$ , one has

$$\begin{split} & \left[ \mathcal{M}_{1} \left( \mathbf{X} \times_{2} U_{2}^{\top} \times_{3} U_{3}^{\top} \right) \right]_{i_{1},(i_{2}-1)r_{3}+i_{3}} = \left( \mathbf{X} \times_{2} U_{2}^{\top} \times_{3} U_{3}^{\top} \right)_{i_{1},i_{2},i_{3}} \\ &= \sum_{j_{2}=1}^{p_{2}} \sum_{j_{3}=1}^{p_{3}} X_{i_{1},j_{2},j_{3}} (U_{2})_{j_{2},i_{2}} (U_{3})_{j_{3},i_{3}} \\ &= \sum_{j_{2}=1}^{p_{2}} \sum_{j_{3}=1}^{p_{3}} \left( \mathcal{M}_{1}(\mathbf{X}) \right)_{i_{1},(j_{2}-1)p_{3}+j_{3}} \cdot (U_{2} \otimes U_{3})_{(j_{2}-1)p_{3}+j_{3},(i_{2}-1)r_{3}+i_{3}} \\ &= \left( \mathcal{M}_{1}(\mathbf{X}) \cdot (U_{2} \otimes U_{3}) \right)_{i_{1},(i_{2}-1)r_{3}+i_{3}} \cdot \end{split}$$

This shows (76).

• The proof for (77) is essentially the same as (76) as we only need to check each entries of the terms in (77) are equal. For (78), let

$$U_2 = \sum_i \sigma_{2i} \cdot \alpha_{2i} \beta_{2i}^{\mathsf{T}}, \quad U_3 = \sum_j \sigma_{3j} \cdot \alpha_{3j} \beta_{3j}^{\mathsf{T}}$$

be the singular value decompositions. Then it is not hard to see the singular value decomposition of  $U_2 \otimes U_3$  can be written as

$$U_2 \otimes U_3 = \sum_{i,j} \sigma_i \sigma_j \cdot (\alpha_{i2} \otimes \alpha_{j3}) (\beta_{i2} \otimes \beta_{j3})^\top,$$

so that the singular values of  $U_{\otimes}U_3$  are  $\{\sigma_i \cdot \sigma_j\}$ . Then

$$||U_2 \otimes U_3|| = \max_{i,j} \sigma_i \sigma_j = \left(\max_i \sigma_i\right) \cdot \left(\max_j \sigma_j\right) = ||U_2|| \cdot ||U_3||,$$

$$\|U_2 \otimes U_3\|_{\mathrm{F}}^2 = \sum_{i,j} \sigma_i^2 \sigma_j^2 = \left(\sum_i \sigma_i^2\right) \cdot \left(\sum_j \sigma_j^2\right) = \|U_2\|_{\mathrm{F}}^2 \cdot \|U_3\|_{\mathrm{F}}^2,$$

and

$$\sigma_{\min}(U_2 \otimes U_3) = \min_{i,j} \sigma_i \sigma_j = \left(\min_i \sigma_i\right) \cdot \left(\min_j \sigma_j\right) = \sigma_{\min}(U_2) \cdot \sigma_{\min}(U_3).$$

$$P_{I_{p_2} \otimes U_3} + P_{I_{p_2} \otimes U_{3\perp}} = (I_{p_2} \otimes U_3)(I_{p_2} \otimes U_3)^{\top} + (I_{p_2} \otimes U_{3\perp})(I_{p_2} \otimes U_{3\perp})^{\top}$$

$$\stackrel{(77)}{=} (I_{p_2} I_{p_2}^{\top}) \otimes (U_3 U_3^{\top}) + (I_{p_2} I_{p_2}^{\top}) \otimes (U_3 \perp U_{3\perp}^{\top}) = I_{p_2} \otimes (U_3 U_3^{\top} + U_3 \perp U_{3\perp}^{\top})$$

$$= I_{p_2} \otimes I_{p_3} = I_{p_2 p_3}.$$

The other identity can be shown similarly.  $\Box$ 

#### C.4 Proof of Lemma 5

The key idea for the proof of this lemma is via  $\varepsilon$ -net. By Lemma 7, for k = 1, 2, 3, there exist  $\varepsilon$ -nets:  $V_k^{(1)}, \dots, V_k^{(N_k)}$  for  $\{V_k \in \mathbb{R}^{p_k \times r_k} : ||V_k|| \le 1\}, |N_k| \le ((4+\varepsilon)/\varepsilon)^{p_k r_k}$ , such that

For any  $V \in \mathbb{R}^{p_k \times r_k}$  satisfying  $||V|| \le 1$ , there exists  $V_k^{(j)}$  such that  $||V_k^{(j)} - V|| \le \varepsilon$ .

For fixed  $V_2^{(i)}$  and  $V_3^{(j)}$ , we consider

$$Z_1^{(ij)} = \mathcal{M}_1 \left( \mathbf{Z} \times_2 (V_2^{(i)})^\top \times_3 (V_3^{(j)})^\top \right) \in \mathbb{R}^{p_1 \times (r_2 r_3)}.$$

Clearly, each row of  $Z_1^{(ij)}$  follows a joint Gaussian distribution:  $N\left(0, \left(V_2^{(i)^{\intercal}}V_2^{(i)}\right) \otimes \left(V_3^{(j)^{\intercal}}V_3^{(j)}\right)\right)$ , and  $\left\|\left(V_2^{(i)^{\intercal}}V_2^{(i)}\right) \otimes \left(V_3^{(j)^{\intercal}}V_3^{(j)}\right)\right\| \leq 1$ . Then by random matrix theory (e.g. [62]),

$$P\left(\|Z_1^{(ij)}\| \le \sqrt{p_1} + \sqrt{r_2 r_3} + t\right) \ge 1 - 2\exp(-t^2/2).$$

Then we further have

$$P\left(\max_{i,j} \|Z_1^{(ij)}\| \le \sqrt{p_1} + \sqrt{r_2 r_3} + x\right) \ge 1 - 2((4+\varepsilon)/\varepsilon)^{p_2 r_2 + p_3 r_3} \exp(-x^2/2), \tag{82}$$

for all x > 0. Now, we assume

$$V_{2}^{*}, V_{3}^{*} = \underset{V_{2} \in \mathbb{R}^{p_{2} \times r_{2}}, V_{3} \in \mathbb{R}^{p_{3} \times r_{3}}}{\arg \max} \left\| \mathcal{M}_{1} \left( \mathbf{Z} \times_{2} V_{2}^{\top} \times_{3} V_{3}^{\top} \right) \right\|,$$

$$M = \underset{V_{2} \in \mathbb{R}^{p_{2} \times r_{2}}, V_{3} \in \mathbb{R}^{p_{3} \times r_{3}}}{\max} \left\| \mathcal{M}_{1} \left( \mathbf{Z} \times_{2} V_{2}^{\top} \times_{3} V_{3}^{\top} \right) \right\|.$$

$$\|V_{2}\| \leq 1, \|V_{3}\| \leq 1$$

By definition of the  $\varepsilon$ -net, we can find  $1 \le i \le N_2$  and  $1 \le j \le N_3$  such that  $||V_2^{(i)} - V_2^*|| \le \varepsilon$  and  $||V_3^{(i)} - V_3^*|| \le \varepsilon$ . In this case under (82),

$$M = \left\| \mathcal{M}_{1} \left( \mathbf{Z} \times_{2} (V_{2}^{*})^{\top} \times_{3} (V_{3}^{*})^{\top} \right) \right\|$$

$$\leq \left\| \mathcal{M}_{1} \left( \mathbf{Z} \times_{2} (V_{2}^{(i)})^{\top} \times_{3} (V_{3}^{(j)})^{\top} \right) \right\| + \left\| \mathcal{M}_{1} \left( \mathbf{Z} \times_{2} (V^{*} - V_{2}^{(i)})^{\top} \times_{3} (V_{3}^{(j)})^{\top} \right) \right\|$$

$$+ \left\| \mathcal{M}_{1} \left( \mathbf{Z} \times_{2} (V_{2}^{*})^{\top} \times_{3} (V_{3}^{*} - V_{3}^{(j)})^{\top} \right) \right\|$$

$$\leq \sqrt{p_{1}} + \sqrt{r_{2}r_{3}} + x + \varepsilon M + \varepsilon M,$$

Therefore, we have

$$P\left(M \le \frac{\sqrt{p_1} + \sqrt{r_2 r_3} + x}{1 - 2\varepsilon}\right) \ge 1 - 2((4 + \varepsilon)/\varepsilon)^{p_2 r_2 + p_3 r_3} \exp(-x^2/2).$$

By setting  $\varepsilon = 1/3$ ,  $x^2 = 2\log(13)(p_2r_2 + p_3r_3)(1+t)$  for some large constant C > 0, we have proved the first part of the lemma.

The proof for the second part is similar. For any given  $V_k \in \mathbb{R}^{p_k \times r_k}$  satisfying  $||V_1||, ||V_2||, ||V_3|| \le 1$ , we have  $||V_1 \otimes V_2 \otimes V_3|| \le 1$ . By Lemma 8, we know

$$P\left(\left\|\mathbf{Z} \times_{1} V_{1}^{\top} \times_{2} V_{2}^{\top} \times_{3} V_{3}^{\top}\right\|_{F}^{2} - \|V_{1} \otimes V_{2} \otimes V_{3}\|_{F}^{2}\right)$$

$$\geq 2\sqrt{t \|(V_{1}^{\top} V_{1}) \otimes (V_{2}^{\top} V_{2}) \otimes (V_{3}^{\top} V_{3})} + 2t \|V_{1} \otimes V_{2} \otimes V_{3}\|^{2} \leq \exp(-t).$$

Since  $||V_1 \otimes V_2 \otimes V_3|| \le 1$ ,  $||V_1 \otimes V_2 \otimes V_3||_F^2 = ||V_1||_F^2 ||V_2||_F^2 ||V_3||_F^2 \le r_1 r_2 r_3$ , then

$$\begin{aligned} &\|(V_1^\top V_1) \otimes (V_2^\top V_2) \otimes (V_3^\top V_3)\|_{\mathrm{F}}^2 = \|V_1^\top V_1\|_{\mathrm{F}}^2 \|V_2^\top V_2\|_{\mathrm{F}}^2 \|V_3^\top V_3\|_{\mathrm{F}}^2 \\ &= \left(\sum_{i=1}^{r_1} \sigma_i^4(V_1)\right) \left(\sum_{i=1}^{r_2} \sigma_i^4(V_1)\right) \left(\sum_{i=1}^{r_3} \sigma_i^4(V_1)\right) \leq r_1 r_2 r_3, \end{aligned}$$

we have for any fixed  $V_1, V_2, V_3$  and x > 0 that

$$P\left(\left\|\mathbf{Z} \times_{1} V_{1}^{\top} \times_{2} V_{2}^{\top} \times_{3} V_{3}^{\top}\right\|_{F}^{2} \ge r_{1} r_{2} r_{3} + 2\sqrt{r_{1} r_{2} r_{3} x} + 2x\right) \le \exp(-x).$$

By geometric inequality,  $2\sqrt{r_1r_2r_3x} \le r_1r_2r_3 + x$ , then we further have

$$P\left(\left\|\mathbf{Z} \times_{1} V_{1}^{\top} \times_{2} V_{2}^{\top} \times_{3} V_{3}^{\top}\right\|_{F}^{2} \ge 2r_{1}r_{2}r_{3} + 3x\right) \le \exp(-x).$$

The rest proof for this lemma is similar to the first part. By Lemma 7, one can find three  $\varepsilon$ -nets:  $V_k^{(1)}, \ldots, V_k^{(N_k)}$  for  $\{V_k \in \mathbb{R}^{p_k \times r_k} : ||V_k|| \le 1\}$  such that  $|N_k| \le ((4+2\varepsilon)/\varepsilon)^{p_k r_k}$ , k = 1, 2, 3. Then

by probability union bound,

$$\max_{V_1^{(a)}, V_2^{(b)}, V_3^{(c)}} P\left( \left\| \mathbf{Z} \times_1 V_1^{\top} \times_2 V_2^{\top} \times_3 V_3^{\top} \right\|_{\mathrm{F}}^2 \ge 2r_1 r_2 r_3 + 3x \right) \\
\le \exp(-x) \cdot \left( (4 + \varepsilon)/\varepsilon \right)^{p_1 r_1 + p_2 r_2 + p_3 r_3}.$$
(83)

When the inequality above holds, we suppose

$$(V_1^*, V_2^*, V_3^*) = \underset{\substack{V_k \in \mathbb{R}^{p_k \times r_k} \\ \|V_k\| < 1}}{\operatorname{arg\,max}} \left\| \mathbf{Z} \times_1 V_1^\top \times_2 V_2^\top \times_3 V_3^\top \right\|, \quad \text{and} \quad T = \left\| \mathbf{Z} \times_1 (V_1^*)^\top \times_2 (V_2^*)^\top \times_3 (V_3^*)^\top \right\|.$$

Then we can find  $V_1^{(a)}, V_2^{(b)}, V_3^{(c)}$  in the corresponding  $\varepsilon$ -nets such that

$$||V_1^* - V_1^{(a)}|| \le \varepsilon, \quad ||V_2^* - V_2^{(b)}|| \le \varepsilon, \quad ||V_3^* - V_3^{(c)}|| \le \varepsilon.$$

Then

$$T = \left\| \mathbf{Z} \times_{1} (V_{1}^{*})^{\top} \times_{2} (V_{2}^{*})^{\top} \times_{3} (V_{3}^{*})^{\top} \right\|$$

$$\leq \left\| \mathbf{Z} \times_{1} (V_{1}^{(a)})^{\top} \times_{2} (V_{2}^{(b)})^{\top} \times_{3} (V_{3}^{(c)})^{\top} \right\| + \left\| \mathbf{Z} \times_{1} (V_{1}^{(a)} - V_{1}^{*})^{\top} \times_{2} (V_{2}^{*})^{\top} \times_{3} (V_{3}^{*})^{\top} \right\|$$

$$+ \left\| \mathbf{Z} \times_{1} (V_{1}^{(a)})^{\top} \times_{2} (V_{2}^{(b)} - V_{2}^{*})^{\top} \times_{3} V_{3}^{\top} \right\| + \left\| \mathbf{Z} \times_{1} (V_{1}^{(a)})^{\top} \times_{2} (V_{2}^{(b)})^{\top} \times_{3} (V_{3}^{(c)} - V_{3}^{*})^{\top} \right\|$$

$$\leq 2r_{1}r_{2}r_{3} + 3t + \left( \|V_{1}^{*} - V_{1}^{(a)}\| + \|V_{2}^{*} - V_{2}^{(b)}\| + \|V_{3}^{*} - V_{3}^{(c)}\| \right) \cdot T,$$

which implies  $T \leq (2r_1r_2r_3 + 3x)/(1 - 3\varepsilon)$  provided that  $\varepsilon < 1/3$  and (83) holds. Let  $\varepsilon = 1/9$ ,  $x = (1+t)\log(37) \cdot (p_1r_1 + p_2r_2 + p_3r_3)$  for some large constant C > 0, by (83) again we have

$$\mathbb{P}\left(T \ge Cr_1r_2r_3 + C(1+t)(p_1r_1 + p_2r_2 + p_3r_3)\right) 
\ge \exp(-Ct(p_1r_1 + p_2r_2 + p_3r_3))$$
(84)

for some uniform constant C > 0. thus we have finished the proof for (80).  $\square$ 

#### C.5 Proof of Lemma 6

$$||P_{\hat{U}_2}X|| \le ||P_{\hat{U}_2}(X+Z)|| + ||Z|| = \sigma_{r+1}(Y) + ||Z|| = \min_{\substack{\tilde{X} \in \mathbb{R}^{p_1 \times p_2} \\ \operatorname{rank}(\tilde{X}) \le r}} ||Y - \tilde{X}|| + ||Z||$$

$$\leq ||Y - X|| + ||Z|| = 2||Z||.$$

Since rank  $(P_{U_2}X) \leq \operatorname{rank}(X) \leq r$ , it is clear that

$$||P_{U_2}X||_{\mathcal{F}} \le 2\sqrt{r}||Z||;$$

meanwhile,

$$\begin{aligned} \left\| P_{\hat{U}_{2}} X \right\|_{\mathcal{F}} &\leq \left\| P_{\hat{U}_{2}} (X+Z) \right\|_{\mathcal{F}} + \|Z\|_{\mathcal{F}} = \left( \sum_{i=r+1}^{p_{1} \wedge p_{2}} \sigma_{i}^{2}(Y) \right)^{1/2} + \|Z\|_{\mathcal{F}} \\ &\leq \min_{\substack{\tilde{X} \in \mathbb{R}^{p_{1} \times p_{2}} \\ \operatorname{rank}(\tilde{X}) \leq r}} \|Y - \tilde{X}\|_{\mathcal{F}} + \|Z\|_{\mathcal{F}} \leq \|Y - X\|_{\mathcal{F}} + \|Z\|_{\mathcal{F}} \leq 2\|Z\|_{\mathcal{F}}, \end{aligned}$$

which has proved this lemma.  $\Box$ 

#### C.6 Proof of Lemma 7

• We first consider the  $\varepsilon$ -net for  $\mathcal{X}_{p_1,p_2}$ . Note that  $\mathcal{X}_{p_1,p_2}$  is a convex set in  $\mathbb{R}^{p_1 \times p_2}$ , we sequentially pick matrices from  $\mathcal{X}_{p_1,p_2}$ , say  $X^{(1)}, X^{(2)}, \ldots$  satisfying the following criterion: for each time t, the picked matrix satisfies  $\min_{t' \leq t} \|X^{(t)} - X^{(t-1)}\|_{\bullet} \geq \varepsilon$ , i.e., the distances from  $X^{(t)}$  to all the other selected matrices are at least  $\varepsilon$ . We stop the selection process until it is not possible to select the next matrix satisfying such criterion.

Suppose now  $X^{(1)}, \ldots, X^{(N)}$  are all we have selected. Since it is not possible to select another matrix from  $\mathcal{X}_{p_1,p_2}$  which meets the criterion, all matrices in  $\mathcal{X}_{p_1,p_2}$  must be within  $\varepsilon$  of some selected matrix in  $\{X^{(1)}, \ldots, X^{(N)}\}$ , thus

$$\mathcal{X}_{p_1,p_2} \subseteq \bigcup_{i=1}^N B(X^{(i)},\varepsilon).$$

Here  $B(X^{(i)}, \varepsilon) = \{X \in \mathbb{R}^{p_1 \times p_2} : ||X - X^{(i)}||_{\bullet} \le \varepsilon\}$  is the closed ball with center  $X^{(i)}$  and radius  $\varepsilon$ , Therefore,  $\{X^{(1)}, \dots, X^{(N)}\}$  is a  $\varepsilon$ -net.

On the other hand, for any  $1 \le i < j \le N$ ,  $||X^{(i)} - X^{(j)}||_{\bullet} \ge \varepsilon$ , so

$$\{X \in \mathbb{R}^{p_1 \times p_2} : ||X||_{\bullet} \le 1 + \varepsilon/2\} \supseteq \bigcup_{i=1}^N B(X^{(i)}, \varepsilon/2),$$

and  $B(X^{(i)}, \varepsilon/2) \cap B(X^{(j)}, \varepsilon/2)$  contains at most one matrix for any  $1 \le i < j \le N$ . Therefore,

$$(1 + \varepsilon/2)^{p_1 p_2} \operatorname{vol}(\mathcal{X}_{p_1, p_2}) = \operatorname{vol}(\{X \in \mathbb{R}^{p_1 \times p_2} : ||X||_{\bullet} \le 1 + \varepsilon/2\})$$

$$\leq \sum_{i=1}^{N} \operatorname{vol}(B^{(i)}, \varepsilon/2) = N(\varepsilon/2)^{p_1 p_2} \operatorname{vol}(\mathcal{X}_{p_1, p_2}),$$
(85)

which implies  $N \leq ((2+\varepsilon)/\varepsilon)^{p_1p_2}$ .

• By the first part of this lemma, there exist  $(\varepsilon/2)$ -nets  $\bar{\mathcal{X}}_{p_1,r}$  and  $\bar{\mathcal{X}}_{r,p_2}$  for  $\{\|X \in \mathbb{R}^{p_1 \times r} : \|X\| \le 1\}$  and  $\{\|X \in \mathbb{R}^{r \times p_2} : \|X\| \le 1\}$ , such that

$$\left|\bar{\mathcal{X}}_{p_1,r}\right| \leq \left(\frac{4+\varepsilon}{\varepsilon}\right)^{p_1r}, \quad \left|\bar{\mathcal{X}}_{r,p_2}\right| \leq \left(\frac{4+\varepsilon}{\varepsilon}\right)^{p_2r}.$$

Next, we argue that

$$\bar{\mathcal{F}}_{p_1,p_2,r} := \left\{ X \cdot Y : X \in \bar{\mathcal{X}}_{p_1,r}, Y \in \bar{\mathcal{X}}_{r,p_2} \right\}$$

is an  $\varepsilon$ -net for  $\mathcal{F}_{p_1,p_2,r}$  in the spectral norm. Actually for any  $X \in \mathcal{F}_{p_1,p_2,r}$ , we can find A,B such that  $X=A\cdot B,\ A\in\mathbb{R}^{p_1\times r}, \|A\|\leq 1; B\in\mathbb{R}^{r\times p_2}, \|B\|\leq 1$ . Then we can find  $A^*\in\bar{\mathcal{X}}_{p_1,r}$  and  $B^*\in\bar{\mathcal{X}}_{r,p_2}$  such that  $\|A-A^*\|\leq \varepsilon/2, \|B-B^*\|\leq \varepsilon/2$ , thus  $A^*B^*\in\bar{\mathcal{F}}_{p_1,p_2,r}$  satisfies

$$||X - A^*B^*|| = ||(AB - AB^*) + (AB^* - A^*B^*)||$$
  
$$\leq ||A|| \cdot ||B - B^*|| + ||A - A^*|| \cdot ||B^*|| \leq 1 \cdot \varepsilon/2 + 1 \cdot \varepsilon/2 = \varepsilon.$$

Note that  $|\bar{\mathcal{X}}_{p_1,p_2,r}| \leq |\bar{\mathcal{X}}_{p_1,p_2,r}| \cdot |\bar{\mathcal{X}}_{p_1,p_2,r}| \leq ((4+\varepsilon)/\varepsilon)^{r(p_1+p_2)}$ , this has finished the proof of this lemma.  $\square$ 

#### C.7 Proof of Lemma 8

Suppose  $A = U\Sigma V^{\top}$  is the singular value decomposition of A. Since U, V are orthogonal and  $u \stackrel{iid}{\sim} N(0,1)$ ,  $||Au||_2^2$  has the same distribution as  $\sum_{i=1}^{p\wedge n} \sigma_i(A)^2 u_i^2$ . By the exponential probability for general chi-square distribution (Lemma 1 in [65]), we have

$$P\left(\sum_{i=1}^{p \wedge n} \sigma_{i}^{2}(A)u_{i}^{2} - \sum_{i=1}^{p \wedge n} \sigma_{i}^{2}(A) \leq -2\sqrt{t \sum_{i=1}^{p \wedge n} \sigma_{i}^{4}(A)}\right) \leq \exp(-x);$$

$$P\left(\sum_{i=1}^{p \wedge n} \sigma_{i}^{2}(A)u_{i}^{2} - \sum_{i=1}^{p \wedge n} \sigma_{i}^{2}(A) \geq 2\sqrt{t \sum_{i=1}^{p \wedge n} \sigma_{i}^{4}(A) + 2t \max \sigma_{i}^{2}(A)}\right) \leq \exp(-x),$$

which has finished the proof for Lemma 5 since  $||A^{\top}A||_{\mathrm{F}}^2 = \sum_{i=1}^{p \wedge n} \sigma_i^4(A)$ , and  $||A|| = \max \sigma_i(A)$ .

#### C.8 Proof of Lemma 10

Clearly, it suffices to prove the claim for i=0, i.e., under  $H_0$ . Let  $G=(V,E)\sim H_0$  with  $V=\{1,2,\ldots,N\}$  and **A** denote its adjacency tensor, meaning that there is a clique of size

 $\kappa_N$  planted in the subset  $\{1, 2, \dots, \lfloor N/2 \rfloor\}$ . Recall that N = 3p for an even integer p. The vertices set of the planted clique is denoted by  $C \subset \{1, \dots, \frac{3p}{2}\}$  with  $|C| = \kappa_N = 20k$  where  $k = \lfloor p^{(1-\tau)/2} \rfloor$ . Recall  $V_1, V_2, V_3$  and define

$$C_i := C \cap V_i, \quad j = 1, 2, 3$$

which represents the subsets of clique vertices in  $V_1, V_2, V_3$ . If  $\mathbf{Y} = \mathcal{T}(\mathbf{A}, \mathbf{\Xi}^-, \mathbf{\Xi}^+) \in \mathbb{R}^{p \times p \times p}$ , it is clear that, under  $H_0, \mathbf{X} = \mathbb{E}(\mathbf{Y}|C)$  is a sparse tensor with supports  $S_1(\mathbf{X}) = C_1 \subset [p/2], S_2(\mathbf{X}) = C_2 - \frac{p}{2} \subset [p/2]$  and  $S_3(\mathbf{X}) = C_3 - p \subset [p/2]$ . We show that the sizes of  $S_1(\mathbf{X}), S_2(\mathbf{X}), S_3(\mathbf{X})$  are lower bounded by k with high probability.

**Lemma 11.** There exists an event  $\mathcal{E}$  on which  $\min\{|S_1(\mathbf{X})|, |S_2(\mathbf{X})|, |S_3(\mathbf{X})|\} \geq k$  and

$$\mathbb{P}(\mathcal{E}) \ge 1 - 6k(0.86)^{2.5k}.$$

For any fixed realization  $G \sim H_0$  with set of clique vertices  $C = C_1 \cup C_2 \cup C_3$  (with corresponding supports  $S_k := S_k(\mathbf{X}), k = 1, 2, 3$ ), we generate a Gaussian random tensor  $\tilde{\mathbf{Y}} \in \mathbb{R}^{p \times p \times p}$  with independent entries such that

$$\tilde{Y}(a,b,c) \sim \mathcal{N}(\mu,1)$$
 if  $(a,b,c) \in S_1 \times S_2 \times S_3$ ;  $\tilde{Y}(a,b,c) \sim \mathcal{N}(0,1)$  otherwise,

where  $S_1 = C_1, S_2 = C_2 - \frac{p}{2}$  and  $S_3 = C_3 - p$ . By Lemma 9, we have

$$\mathrm{d_{TV}}\Big(\mathcal{L}\big(Y(a,b,c)\big|C\big),\mathcal{L}\big(\tilde{Y}(a,b,c)\big|C\big)\Big) \leq e^{(1-M^2)/2}, \quad \forall \ a,b,c \in \{1,2,\ldots,p\}.$$

As a result, since  $M = \sqrt{8 \log N}$  and p = N/3,

$$d_{\text{TV}}(\mathcal{L}(\mathbf{Y}), \mathcal{L}(\tilde{\mathbf{Y}})) = \mathbb{E}_C d_{\text{TV}}(\mathcal{L}(\mathbf{Y}|C), \mathcal{L}(\tilde{\mathbf{Y}}|C))$$

$$= \mathbb{E}_C \sum_{i=1}^p d_{\text{TV}}(\mathcal{L}(Y(a, b, c)|C), \mathcal{L}(\tilde{Y}(a, b, c)|C)) \leq p^3 e^{(1-M^2)/2} \leq \frac{\sqrt{e}}{27N}.$$

Now we show that  $\mathcal{L}(\tilde{\mathbf{Y}}|\mathcal{E})$  is a mixture over  $\{\mathbb{P}_{\mathbf{X}}, \mathbf{X} \in \mathcal{M}_0(\boldsymbol{p}, k, \boldsymbol{r}, \lambda)\}$  with  $\lambda = \frac{p^{3(1-\tau)/4}}{2\sqrt{8\log 3p}}$ . Indeed, for any fixed C, let  $\tilde{\mathbf{X}} = \mathbb{E}(\tilde{\mathbf{Y}}|C)$ . Then,

$$\tilde{X}(a,b,c) = \mathbb{E}(\tilde{Y}(a,b,c)|C) = \mu, \quad \forall (a,b,c) \in S_1 \times S_2 \times S_3.$$

Recall that on  $\mathcal{E}$ ,  $\min\{|S_1|, |S_2|, |S_3|\} \geq k$ . Therefore,  $\tilde{\mathbf{X}}$  is of rank 1 and on  $\mathcal{E}$ ,

$$\min \left\{ \sigma_{\min} \left( \mathcal{M}_1(\tilde{\mathbf{X}}) \right), \sigma_{\min} \left( \mathcal{M}_2(\tilde{\mathbf{X}}) \right), \sigma_{\min} \left( \mathcal{M}_3(\tilde{\mathbf{X}}) \right) \right\}$$

$$\geq \mu \sqrt{|S_1||S_2||S_3|} \geq \mu k^{3/2} \geq \mu \lfloor p^{3(1-\tau)/4} \rfloor = \frac{\lfloor p^{3(1-\tau)/4} \rfloor}{2\sqrt{8\log 3p}}$$

since  $\mu = \frac{1}{2M}$ . The above fact indicates that  $\tilde{\mathbf{X}} \in \mathcal{M}_0(\boldsymbol{p}, k, \boldsymbol{r}, \lambda)$ . In other words, under  $H_0$ , for any C conditioned on  $\mathcal{E}$ , there exists  $\mathbf{X}(C) \in \mathcal{M}_0(\boldsymbol{p}, k, \boldsymbol{r}, \lambda)$  such that  $\mathcal{L}(\tilde{\mathbf{Y}}|C) = \mathbb{P}_{\mathbf{X}(C)}$ . Define the probability distribution  $\pi_0 = \mathcal{L}(\mathbf{X}(C)|\mathcal{E})$  supported on  $\mathcal{M}_0(\boldsymbol{p}, k, \boldsymbol{r}, \lambda)$ . Then  $\mathcal{L}(\tilde{\mathbf{Y}}|\mathcal{E}) = \mathbb{P}_{\pi_0}$  and

$$d_{TV}(\mathcal{L}(\mathbf{Y}), \mathbb{P}_{\pi_0}) \leq d_{TV}(\mathcal{L}(\mathbf{Y}), \mathcal{L}(\tilde{\mathbf{Y}})) + TV(\mathcal{L}(\tilde{\mathbf{Y}}), \mathbb{P}_{\pi_0})$$
$$\leq \frac{\sqrt{e}}{27N} + \mathbb{P}(\mathcal{E}^c) \leq \frac{\sqrt{e}}{27N} + 6k(0.86)^{2.5k}.$$

#### C.9 Proof of Lemma 11

Recall that  $\kappa \leq \sqrt{N/2}$  and N=3p. Let  $N_1 := N/2 = \frac{3p}{2}$ . Since C is uniformly chosen from  $\{1, 2, \dots, \frac{3p}{2}\}$  and  $C_1 \subset \{1, 2, \dots, \frac{p}{2}\}$ , we have

$$\mathbb{P}\Big(|C_1| \leq \frac{\kappa}{8}\Big) \leq \frac{\sum_{s=0}^{\kappa/8} \binom{p/2}{s} \binom{p}{\kappa-s}}{\binom{N_1}{\kappa}} \leq \frac{\kappa+1}{8} \frac{\binom{p/2}{\kappa/8} \binom{p}{7\kappa/8}}{\binom{N_1}{\kappa}} \\
= \frac{\kappa+1}{8} \binom{\kappa}{\kappa/8} \frac{(p)(p-1)\dots(p-7\kappa/8+1)(p/2)(p/2-1)\dots(p/2-\kappa/8+1)}{(3p/2)(3p/2-1)\dots(3p/2-\kappa+1)} \\
\leq \frac{\kappa+1}{8} (8e)^{\kappa/8} (\frac{2}{3})^{\kappa} = \frac{\kappa}{8} \Big(8e \cdot 2^8/3^8\Big)^{\kappa/8} \leq \frac{\kappa+1}{8} (0.86)^{\kappa/8}$$

where we used the fact  $\binom{p/2}{s}\binom{p}{\kappa-s}$  increases for  $0 \le s \le \kappa/8$  and inequality  $\binom{n}{k} \le (ne/k)^k$ . Therefore, with probability at least  $1 - \frac{\kappa+1}{4}(0.86)^{\kappa/8}$ ,

$$\frac{\kappa}{8} \le |C_1|, |C_2|, |C_3| \le \frac{7\kappa}{8}.$$

Recall that  $\kappa = 20k$  and we conclude the proof.