

Tensor Robust Principal Component Analysis with A New Tensor Nuclear Norm

Canyi Lu, Jiashi Feng, Yudong Chen, Wei Liu, Member, IEEE, Zhouchen Lin, *Fellow, IEEE*, and Shuicheng Yan, *Fellow, IEEE*

Abstract—In this paper, we consider the Tensor Robust Principal Component Analysis (TRPCA) problem, which aims to exactly recover the low-rank and sparse components from their sum. Our model is based on the recently proposed tensor-tensor product (or t-product) [15]. Induced by the t-product, we first rigorously deduce the tensor spectral norm, tensor nuclear norm, and tensor average rank, and show that the tensor nuclear norm is the convex envelope of the tensor average rank within the unit ball of the tensor spectral norm. These definitions, their relationships and properties are consistent with matrix cases. Equipped with the new tensor nuclear norm, we then solve the TRPCA problem by solving a convex program and provide the theoretical guarantee for the exact recovery. Our TRPCA model and recovery guarantee include matrix RPCA as a special case. Numerical experiments verify our results, and the applications to image recovery and background modeling problems demonstrate the effectiveness of our method.

Index Terms—Tensor robust PCA, convex optimization, tensor nuclear norm, tensor singular value decomposition

1 INTRODUCTION

PRINCIPAL Component Analysis (PCA) is a fundamental approach for data analysis. It exploits low-dimensional structure in high-dimensional data, which commonly exists in different types of data, *e.g.*, image, text, video and bioinformatics. It is computationally efficient and powerful for data instances which are mildly corrupted by small noises. However, a major issue of PCA is that it is brittle to be grossly corrupted or outlying observations, which are ubiquitous in real-world data. To date, a number of robust versions of PCA have been proposed, but many of them suffer from a high computational cost.

The Robust PCA [3] is the first polynomial-time algorithm with strong recovery guarantees. Suppose that we are given an observed matrix $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$, which can be decomposed as $\mathbf{X} = \mathbf{L}_0 + \mathbf{E}_0$, where \mathbf{L}_0 is low-rank and \mathbf{E}_0 is sparse. It is shown in [3] that if the singular vectors of \mathbf{L}_0 satisfy some incoherent conditions, *e.g.*, \mathbf{L}_0 is low-rank and \mathbf{E}_0 is sufficiently sparse, then \mathbf{L}_0 and \mathbf{E}_0 can be exactly recovered with high probability by solving the following convex problem

$$\min_{\mathbf{L}, \mathbf{E}} \|\mathbf{L}\|_* + \lambda \|\mathbf{E}\|_1, \text{ s.t. } \mathbf{X} = \mathbf{L} + \mathbf{E}, \quad (1)$$

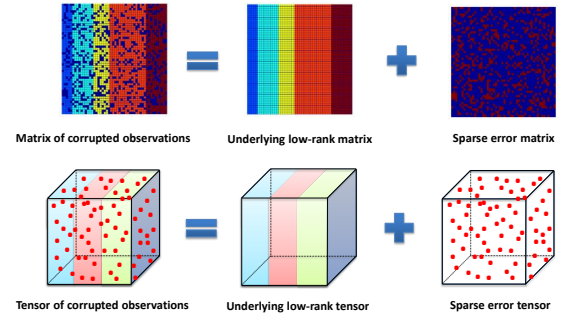


Fig. 1: Illustrations of RPCA [3] (up row) and our Tensor RPCA (bottom row). RPCA: low-rank and sparse matrix decomposition from noisy matrix observations. Tensor RPCA: low-rank and sparse tensor decomposition from noisy tensor observations.

where $\|\mathbf{L}\|_*$ denotes the nuclear norm (sum of the singular values of \mathbf{L}), and $\|\mathbf{E}\|_1$ denotes the ℓ_1 -norm (sum of the absolute values of all the entries in \mathbf{E}). Theoretically, RPCA is guaranteed to work even if the rank of \mathbf{L}_0 grows almost linearly in the dimension of the matrix, and the errors in \mathbf{E}_0 are up to a constant fraction of all entries. The parameter λ is suggested to be set as $1/\sqrt{\max(n_1, n_2)}$ which works well in practice. Algorithmically, program (1) can be solved by efficient algorithms, at a cost not too much higher than PCA. RPCA and its extensions have been successfully applied to background modeling [3], subspace clustering [17], video compressive sensing [31], *etc.*

One major shortcoming of RPCA is that it can only handle 2-way (matrix) data. However, real data is usually multi-dimensional in nature—the information is stored in multi-way arrays known as tensors [16]. For example, a color image is a 3-way object with column, row and color modes; a greyscale video

- C. Lu is with the Department of Electrical and Computer Engineering, Carnegie Mellon University (e-mail: canyilu@gmail.com).
- J. Feng and S. Yan are with the Department of Electrical and Computer Engineering, National University of Singapore, Singapore (e-mail: elefjia@nus.edu.sg; eleyans@nus.edu.sg).
- Y. Chen is with the School of Operations Research and Information Engineering, Cornell University (e-mail: yudong.chen@cornell.edu).
- W. Liu is with the Tencent AI Lab, Shenzhen, China (e-mail: wl2223@columbia.edu).
- Z. Lin is with the Key Laboratory of Machine Perception (MOE), School of Electronics Engineering and Computer Science, Peking University, Beijing 100871, China (e-mail: zlin@pku.edu.cn).

is indexed by two spatial variables and one temporal variable. To use RPCA, one has to first restructure the multi-way data into a matrix. Such a preprocessing usually leads to an information loss and would cause a performance degradation. To alleviate this issue, it is natural to consider extending RPCA to manipulate the tensor data by taking advantage of its multi-dimensional structure.

In this work, we are interested in the Tensor Robust Principal Component (TRPCA) model which aims to exactly recover a low-rank tensor corrupted by sparse errors. See Figure 1 for an intuitive illustration. More specifically, suppose that we are given a data tensor \mathcal{X} , and know that it can be decomposed as

$$\mathcal{X} = \mathcal{L}_0 + \mathcal{E}_0, \quad (2)$$

where \mathcal{L}_0 is low-rank and \mathcal{E}_0 is sparse, and both components are of arbitrary magnitudes. Note that we do not know the locations of the nonzero elements of \mathcal{E}_0 , not even how many there are. Now we consider a similar problem to RPCA. Can we recover the low-rank and sparse components exactly and efficiently from \mathcal{X} ? This is the problem of tensor RPCA studied in this work.

The tensor extension of RPCA is not easy since the numerical algebra of tensors is fraught with hardness results [11], [5], [8]. A main issue is that the tensor rank is not well defined with a tight convex relaxation. Several tensor rank definitions and their convex relaxations have been proposed but each has its limitation. For example, the CP rank [16], defined as the smallest number of rank one tensor decomposition, is generally NP-hard to compute. Also its convex relaxation is intractable. This makes the low CP rank tensor recovery challenging. The tractable Tucker rank [16] and its convex relaxation are more widely used. For a k -way tensor \mathcal{X} , the Tucker rank is a vector defined as $\text{rank}_{\text{tc}}(\mathcal{X}) := (\text{rank}(\mathbf{X}^{\{1\}}), \text{rank}(\mathbf{X}^{\{2\}}), \dots, \text{rank}(\mathbf{X}^{\{k\}}))$, where $\mathbf{X}^{\{i\}}$ is the mode- i matricization of \mathcal{X} [16]. Motivated by the fact that the nuclear norm is the convex envelope of the matrix rank within the unit ball of the spectral norm, the Sum of Nuclear Norms (SNN) [18], defined as $\sum_i \|\mathbf{X}^{\{i\}}\|_*$, is used as a convex surrogate of $\sum_i \text{rank}(\mathbf{X}^{\{i\}})$. Then the work [24] considers the Low-Rank Tensor Completion (LRTC) model based on SNN:

$$\min_{\mathcal{X}} \sum_{i=1}^k \lambda_i \|\mathbf{X}^{\{i\}}\|_*, \text{ s.t. } \mathcal{P}_{\Omega}(\mathcal{X}) = \mathcal{P}_{\Omega}(\mathcal{M}), \quad (3)$$

where $\lambda_i > 0$, and $\mathcal{P}_{\Omega}(\mathcal{X})$ denotes the projection of \mathcal{X} on the observed set Ω . The effectiveness of this approach for image processing has been well studied in [18], [28]. However, SNN is not the convex envelope of $\sum_i \text{rank}(\mathbf{X}^{\{i\}})$ [26]. Actually, the above model can be substantially suboptimal [24]: reliably recovering a k -way tensor of length n and Tucker rank (r, r, \dots, r) from Gaussian measurements requires $O(rn^{k-1})$ observations. In contrast, a certain (intractable) nonconvex formulation needs only $O(rK + nrK)$ observations. A better (but still suboptimal) convexification based on a more balanced matricization is proposed in [24]. The work [13] presents the recovery guarantee for the SNN based tensor RPCA model

$$\min_{\mathcal{L}, \mathcal{E}} \sum_{i=1}^k \lambda_i \|\mathbf{L}^{\{i\}}\|_* + \|\mathcal{E}\|_1, \text{ s.t. } \mathcal{X} = \mathcal{L} + \mathcal{E}. \quad (4)$$

A robust tensor CP decomposition problem is studied in [6]. Though the recovery is guaranteed, the algorithm is nonconvex.

The limitations of existing works motivate us to consider an interesting problem: is it possible to define a new tensor nuclear norm such that it is a tight convex surrogate of certain tensor rank, and thus its resulting tensor RPCA enjoys a similar tight recovery guarantee to that of the matrix RPCA? This work will provide a positive answer to this question. Our solution is inspired by the recently proposed tensor-tensor product (t-product) [15] which is a generalization of the matrix-matrix product. It enjoys several similar properties to the matrix-matrix product. For example, based on t-product, any tensors have the tensor Singular Value Decomposition (t-SVD) and this motivates a new tensor rank, *i.e.*, tensor tubal rank [14]. To recover a tensor of low tubal rank, we propose a new tensor nuclear norm which is rigorously induced by the t-product. First, the tensor spectral norm can be induced by the operator norm when treating the t-product as an operator. Then the tensor nuclear norm is defined as the dual norm of the tensor spectral norm. We further propose the tensor average rank (which is closely related to the tensor tubal rank), and prove that its convex envelope is the tensor nuclear norm within the unit ball of the tensor spectral norm. It is interesting that this framework, including the new tensor concepts and their relationships, is consistent with the one for the matrix cases. Equipped with these new tools, we then study the TRPCA problem which aims to recover the low tubal rank component \mathcal{L}_0 and sparse component \mathcal{E}_0 from noisy observations $\mathcal{X} = \mathcal{L}_0 + \mathcal{E}_0 \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ (this work focuses on the 3-way tensor) by convex optimization

$$\min_{\mathcal{L}, \mathcal{E}} \|\mathcal{L}\|_* + \lambda \|\mathcal{E}\|_1, \text{ s.t. } \mathcal{X} = \mathcal{L} + \mathcal{E}, \quad (5)$$

where $\|\mathcal{L}\|_*$ is our new tensor nuclear norm (see the definition in Section 3). We prove that under certain incoherence conditions, the solution to (5) perfectly recovers the low-rank and the sparse components, provided of course that the tubal rank of \mathcal{L}_0 is not too large, and that \mathcal{E}_0 is reasonably sparse. A remarkable fact, like in RPCA, is that (5) has no tuning parameter either. Our analysis shows that $\lambda = 1/\sqrt{\max(n_1, n_2)n_3}$ guarantees the exact recovery when \mathcal{L}_0 and \mathcal{E}_0 satisfy certain assumptions. As a special case, if \mathcal{X} reduces to a matrix ($n_3 = 1$ in this case), all the new tensor concepts reduce to the matrix cases. Our TRPCA model (5) reduces to RPCA in (1), and also our recovery guarantee in Theorem 4.1 reduces to Theorem 1.1 in [3]. Another advantage of (5) is that it can be solved by polynomial-time algorithms.

The contributions of this work are summarized as follows:

1. Motivated by the t-product [15] which is a natural generalization of the matrix-matrix product, we rigorously deduce a new tensor nuclear norm and some other related tensor concepts, and they own the same relationship as the matrix cases. This is the foundation for the extensions of the models, optimization method and theoretical analyzing techniques from matrix cases to tensor cases.
2. Equipped with the tensor nuclear norm, we theoretically show that under certain incoherence conditions, the solution to the convex TRPCA model (5) perfectly recovers the underlying

low-rank component \mathcal{L}_0 and sparse component \mathcal{E}_0 . RPCA [3] and its recovery guarantee fall into our special cases.

3. We give a new rigorous proof of t-SVD factorization and a more efficient way than [19] for solving TRPCA. We further perform several simulations to corroborate our theoretical results. Numerical experiments on images and videos also show the superiority of TRPCA over RPCA and SNN.

The rest of this paper is structured as follows. Section A gives some notations and preliminaries. Section 3 presents the way for defining the tensor nuclear norm induced by the t-product. Section 4 provides the recovery guarantee of TRPCA and the optimization details. Section 5 presents numerical experiments conducted on synthetic and real data. We conclude this work in Section 6.

2 NOTATIONS AND PRELIMINARIES

2.1 Notations

In this paper, we denote tensors by boldface Euler script letters, *e.g.*, \mathcal{A} . Matrices are denoted by boldface capital letters, *e.g.*, \mathbf{A} ; vectors are denoted by boldface lowercase letters, *e.g.*, \mathbf{a} , and scalars are denoted by lowercase letters, *e.g.*, a . We denote \mathbf{I}_n as the $n \times n$ identity matrix. The fields of real numbers and complex numbers are denoted as \mathbb{R} and \mathbb{C} , respectively. For a 3-way tensor $\mathcal{A} \in \mathbb{C}^{n_1 \times n_2 \times n_3}$, we denote its (i, j, k) -th entry as \mathcal{A}_{ijk} or a_{ijk} and use the Matlab notation $\mathcal{A}(i, :, :)$, $\mathcal{A}(:, i, :)$ and $\mathcal{A}(:, :, i)$ to denote respectively the i -th horizontal, lateral and frontal slice (see definitions in [16]). More often, the frontal slice $\mathcal{A}(:, :, i)$ is denoted compactly as $\mathcal{A}^{(i)}$. The tube is denoted as $\mathcal{A}(i, j, :)$. The inner product between \mathcal{A} and \mathcal{B} in $\mathbb{C}^{n_1 \times n_2}$ is defined as $\langle \mathcal{A}, \mathcal{B} \rangle = \text{Tr}(\mathcal{A}^* \mathcal{B})$, where \mathcal{A}^* denotes the conjugate transpose of \mathcal{A} and $\text{Tr}(\cdot)$ denotes the matrix trace. The inner product between \mathcal{A} and \mathcal{B} in $\mathbb{C}^{n_1 \times n_2 \times n_3}$ is defined as $\langle \mathcal{A}, \mathcal{B} \rangle = \sum_{i=1}^{n_3} \langle \mathcal{A}^{(i)}, \mathcal{B}^{(i)} \rangle$. For any $\mathcal{A} \in \mathbb{C}^{n_1 \times n_2 \times n_3}$, the complex conjugate of \mathcal{A} is denoted as $\text{conj}(\mathcal{A})$ which takes the complex conjugate of each entry of \mathcal{A} . We denote $\lfloor t \rfloor$ as the nearest integer less than or equal to t and $\lceil t \rceil$ as the one greater than or equal to t .

Some norms of vector, matrix and tensor are used. We denote the ℓ_1 -norm as $\|\mathcal{A}\|_1 = \sum_{ijk} |a_{ijk}|$, the infinity norm as $\|\mathcal{A}\|_\infty = \max_{ijk} |a_{ijk}|$ and the Frobenius norm as $\|\mathcal{A}\|_F = \sqrt{\sum_{ijk} |a_{ijk}|^2}$, respectively. The above norms reduce to the vector or matrix norms if \mathcal{A} is a vector or a matrix. For $\mathbf{v} \in \mathbb{C}^n$, the ℓ_2 -norm is $\|\mathbf{v}\|_2 = \sqrt{\sum_i |v_i|^2}$. The spectral norm of a matrix \mathbf{A} is denoted as $\|\mathbf{A}\| = \max_i \sigma_i(\mathbf{A})$, where $\sigma_i(\mathbf{A})$'s are the singular values of \mathbf{A} . The matrix nuclear norm is $\|\mathbf{A}\|_* = \sum_i \sigma_i(\mathbf{A})$.

2.2 Discrete Fourier Transformation

The Discrete Fourier Transformation (DFT) plays a core role in tensor-tensor product introduced later. We give some related background knowledge and notations here. The DFT on $\mathbf{v} \in \mathbb{R}^n$, denoted as $\bar{\mathbf{v}}$, is given by

$$\bar{\mathbf{v}} = \mathbf{F}_n \mathbf{v} \in \mathbb{C}^n, \quad (6)$$

where \mathbf{F}_n is the DFT matrix defined as

$$\mathbf{F}_n = \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & \omega & \omega^2 & \cdots & \omega^{n-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \omega^{n-1} & \omega^{2(n-1)} & \cdots & \omega^{(n-1)(n-1)} \end{bmatrix} \in \mathbb{C}^{n \times n},$$

where $\omega = e^{-\frac{2\pi i}{n}}$ is a primitive n -th root of unity in which $i = \sqrt{-1}$. Note that \mathbf{F}_n/\sqrt{n} is a unitary matrix, *i.e.*,

$$\mathbf{F}_n^* \mathbf{F}_n = \mathbf{F}_n \mathbf{F}_n^* = n \mathbf{I}_n. \quad (7)$$

Thus $\mathbf{F}_n^{-1} = \mathbf{F}_n^*/n$. The above property will be frequently used in this paper. Computing $\bar{\mathbf{v}}$ by using (6) costs $O(n^2)$. A more widely used method is the Fast Fourier Transform (FFT) which costs $O(n \log n)$. By using the Matlab command `fft`, we have $\bar{\mathbf{v}} = \text{fft}(\mathbf{v})$. Denote the circulant matrix of \mathbf{v} as

$$\text{circ}(\mathbf{v}) = \begin{bmatrix} v_1 & v_n & \cdots & v_2 \\ v_2 & v_1 & \cdots & v_3 \\ \vdots & \vdots & \ddots & \vdots \\ v_n & v_{n-1} & \cdots & v_1 \end{bmatrix} \in \mathbb{R}^{n \times n}.$$

It is known that it can be diagonalized by the DFT matrix, *i.e.*,

$$\mathbf{F}_n \cdot \text{circ}(\mathbf{v}) \cdot \mathbf{F}_n^{-1} = \text{Diag}(\bar{\mathbf{v}}), \quad (8)$$

where $\text{Diag}(\bar{\mathbf{v}})$ denotes a diagonal matrix with its i -th diagonal entry as \bar{v}_i . The above equation implies that the columns of \mathbf{F}_n are the eigenvectors of $(\text{circ}(\mathbf{v}))^T$ and \bar{v}_i 's are the corresponding eigenvalues.

Lemma 2.1. [25] *Given any real vector $\mathbf{v} \in \mathbb{R}^n$, the associated $\bar{\mathbf{v}}$ satisfies*

$$\bar{v}_1 \in \mathbb{R} \text{ and } \text{conj}(\bar{v}_i) = \bar{v}_{n-i+2}, \quad i = 2, \dots, \left\lfloor \frac{n+1}{2} \right\rfloor. \quad (9)$$

Conversely, for any given complex $\bar{\mathbf{v}} \in \mathbb{C}^n$ satisfying (9), there exists a real block circulant matrix $\text{circ}(\mathbf{v})$ such that (8) holds.

As will be seen later, the above properties are useful for efficient computation and important for proofs. Now we consider the DFT on tensors. For $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, we denote $\bar{\mathcal{A}} \in \mathbb{C}^{n_1 \times n_2 \times n_3}$ as the result of DFT on \mathcal{A} along the 3-rd dimension, *i.e.*, performing the DFT on all the tubes of \mathcal{A} . By using the Matlab command `fft`, we have

$$\bar{\mathcal{A}} = \text{fft}(\mathcal{A}, [], 3).$$

In a similar fashion, we can compute \mathcal{A} from $\bar{\mathcal{A}}$ using the inverse FFT, *i.e.*,

$$\mathcal{A} = \text{ifft}(\bar{\mathcal{A}}, [], 3).$$

In particular, we denote $\bar{\mathcal{A}} \in \mathbb{C}^{n_1 n_3 \times n_2 n_3}$ as a block diagonal matrix with its i -th block on the diagonal as the i -th frontal slice $\bar{\mathcal{A}}^{(i)}$ of $\bar{\mathcal{A}}$, *i.e.*,

$$\bar{\mathcal{A}} = \text{bdiag}(\bar{\mathcal{A}}) = \begin{bmatrix} \bar{\mathcal{A}}^{(1)} & & & \\ & \bar{\mathcal{A}}^{(2)} & & \\ & & \ddots & \\ & & & \bar{\mathcal{A}}^{(n_3)} \end{bmatrix},$$

where bdiag is an operator which maps the tensor $\bar{\mathcal{A}}$ to the block diagonal matrix $\bar{\mathbf{A}}$. Also, we define the block circulant matrix $\text{bcirc}(\mathcal{A}) \in \mathbb{R}^{n_1 n_3 \times n_2 n_3}$ of \mathcal{A} as

$$\text{bcirc}(\mathcal{A}) = \begin{bmatrix} \mathbf{A}^{(1)} & \mathbf{A}^{(n_3)} & \dots & \mathbf{A}^{(2)} \\ \mathbf{A}^{(2)} & \mathbf{A}^{(1)} & \dots & \mathbf{A}^{(3)} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{A}^{(n_3)} & \mathbf{A}^{(n_3-1)} & \dots & \mathbf{A}^{(1)} \end{bmatrix}.$$

Just like the circulant matrix which can be diagonalized by DFT, the block circulant matrix can be block diagonalized, *i.e.*,

$$(\mathbf{F}_{n_3} \otimes \mathbf{I}_{n_1}) \cdot \text{bcirc}(\mathcal{A}) \cdot (\mathbf{F}_{n_3}^{-1} \otimes \mathbf{I}_{n_2}) = \bar{\mathbf{A}}, \quad (10)$$

where \otimes denotes the Kronecker product and $(\mathbf{F}_{n_3} \otimes \mathbf{I}_{n_1})/\sqrt{n_3}$ is unitary. By using Lemma 2.1, we have

$$\begin{cases} \bar{\mathbf{A}}^{(1)} \in \mathbb{R}^{n_1 \times n_2}, \\ \text{conj}(\bar{\mathbf{A}}^{(i)}) = \bar{\mathbf{A}}^{(n_3-i+2)}, i = 2, \dots, \lfloor \frac{n_3+1}{2} \rfloor. \end{cases} \quad (11)$$

Conversely, for any given $\bar{\mathbf{A}} \in \mathbb{C}^{n_1 \times n_2 \times n_3}$ satisfying (11), there exists a real tensor $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ such that (10) holds. Also, by using (7), we have the following properties which will be used frequently:

$$\|\mathcal{A}\|_F = \frac{1}{\sqrt{n_3}} \|\bar{\mathbf{A}}\|_F, \quad (12)$$

$$\langle \mathcal{A}, \mathcal{B} \rangle = \frac{1}{n_3} \langle \bar{\mathbf{A}}, \bar{\mathbf{B}} \rangle. \quad (13)$$

2.3 T-product and T-SVD

For $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, we define

$$\text{unfold}(\mathcal{A}) = \begin{bmatrix} \mathbf{A}^{(1)} \\ \mathbf{A}^{(2)} \\ \vdots \\ \mathbf{A}^{(n_3)} \end{bmatrix}, \quad \text{fold}(\text{unfold}(\mathcal{A})) = \mathcal{A},$$

where the unfold operator maps \mathcal{A} to a matrix of size $n_1 n_3 \times n_2$ and fold is its inverse operator.

Definition 2.1. (T-product) [15] Let $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ and $\mathcal{B} \in \mathbb{R}^{n_2 \times l \times n_3}$. Then the t-product $\mathcal{A} * \mathcal{B}$ is defined to be a tensor of size $n_1 \times l \times n_3$,

$$\mathcal{A} * \mathcal{B} = \text{fold}(\text{bcirc}(\mathcal{A}) \cdot \text{unfold}(\mathcal{B})). \quad (14)$$

The t-product can be understood from two perspectives. First, in the original domain, a 3-way tensor of size $n_1 \times n_2 \times n_3$ can be regarded as an $n_1 \times n_2$ matrix with each entry being a tube that lies in the third dimension. Thus, the t-product is analogous to the matrix multiplication except that the circular convolution replaces the multiplication operation between the elements. Note that the t-product reduces to the standard matrix multiplication when $n_3 = 1$. This is a key observation which makes our tensor RPCA model shown later involve the matrix RPCA as a special case. Second, the t-product is equivalent to the matrix

Algorithm 1 Tensor-Tensor Product

Input: $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, $\mathcal{B} \in \mathbb{R}^{n_2 \times l \times n_3}$.

Output: $\mathcal{C} = \mathcal{A} * \mathcal{B} \in \mathbb{R}^{n_1 \times l \times n_3}$.

1. Compute $\bar{\mathcal{A}} = \text{fft}(\mathcal{A}, [], 3)$ and $\bar{\mathcal{B}} = \text{fft}(\mathcal{B}, [], 3)$.
2. Compute each frontal slice of $\bar{\mathcal{C}}$ by

$$\bar{\mathcal{C}}^{(i)} = \begin{cases} \bar{\mathbf{A}}^{(i)} \bar{\mathbf{B}}^{(i)}, & i = 1, \dots, \lfloor \frac{n_3+1}{2} \rfloor, \\ \text{conj}(\bar{\mathcal{C}}^{(n_3-i+2)}), & i = \lfloor \frac{n_3+1}{2} \rfloor + 1, \dots, n_3. \end{cases}$$

3. Compute $\mathcal{C} = \text{ifft}(\bar{\mathcal{C}}, [], 3)$.

multiplication in the Fourier domain; that is, $\mathcal{C} = \mathcal{A} * \mathcal{B}$ is equivalent to $\bar{\mathcal{C}} = \bar{\mathcal{A}} \bar{\mathcal{B}}$ due to (10). Indeed, $\mathcal{C} = \mathcal{A} * \mathcal{B}$ implies

$$\begin{aligned} & \text{unfold}(\mathcal{C}) \\ &= \text{bcirc}(\mathcal{A}) \cdot \text{unfold}(\mathcal{B}) \\ &= (\mathbf{F}_{n_3}^{-1} \otimes \mathbf{I}_{n_1}) \cdot ((\mathbf{F}_{n_3} \otimes \mathbf{I}_{n_1}) \cdot \text{bcirc}(\mathcal{A}) \cdot (\mathbf{F}_{n_3}^{-1} \otimes \mathbf{I}_{n_2})) \\ & \quad \cdot ((\mathbf{F}_{n_3} \otimes \mathbf{I}_{n_2}) \cdot \text{unfold}(\mathcal{B})) \\ &= (\mathbf{F}_{n_3}^{-1} \otimes \mathbf{I}_{n_1}) \cdot \bar{\mathbf{A}} \cdot \text{unfold}(\bar{\mathcal{B}}), \end{aligned} \quad (15)$$

where (15) uses (10). Left multiplying both sides with $(\mathbf{F}_{n_3} \otimes \mathbf{I}_{n_1})$ leads to $\text{unfold}(\bar{\mathcal{C}}) = \bar{\mathbf{A}} \cdot \text{unfold}(\bar{\mathcal{B}})$. This is equivalent to $\bar{\mathcal{C}} = \bar{\mathcal{A}} \bar{\mathcal{B}}$. This property suggests an efficient way based on FFT to compute t-product instead of using (14). See Algorithm 1.

The t-product enjoys many similar properties to the matrix-matrix product. For example, the t-product is associative, *i.e.*, $\mathcal{A} * (\mathcal{B} * \mathcal{C}) = (\mathcal{A} * \mathcal{B}) * \mathcal{C}$. We also need some other concepts on tensors extended from the matrix cases.

Definition 2.2. (Conjugate transpose) The conjugate transpose of a tensor $\mathcal{A} \in \mathbb{C}^{n_1 \times n_2 \times n_3}$ is the tensor $\mathcal{A}^* \in \mathbb{C}^{n_2 \times n_1 \times n_3}$ obtained by conjugate transposing each of the frontal slices and then reversing the order of transposed frontal slices 2 through n_3 .

The tensor conjugate transpose extends the tensor transpose [15] for complex tensors. As an example, let $\mathcal{A} \in \mathbb{C}^{n_1 \times n_2 \times 4}$ and its frontal slices be $\mathbf{A}_1, \mathbf{A}_2, \mathbf{A}_3$ and \mathbf{A}_4 . Then

$$\mathcal{A}^* = \text{fold} \left(\begin{bmatrix} \mathbf{A}_1^* \\ \mathbf{A}_4^* \\ \mathbf{A}_3^* \\ \mathbf{A}_2^* \end{bmatrix} \right).$$

Definition 2.3. (Identity tensor) [15] The identity tensor $\mathcal{I} \in \mathbb{R}^{n \times n \times n_3}$ is the tensor with its first frontal slice being the $n \times n$ identity matrix, and other frontal slices being all zeros.

It is clear that $\mathcal{A} * \mathcal{I} = \mathcal{A}$ and $\mathcal{I} * \mathcal{A} = \mathcal{A}$ given the appropriate dimensions. The tensor $\bar{\mathcal{I}} = \text{fft}(\mathcal{I}, [], 3)$ is a tensor with each frontal slice being the identity matrix.

Definition 2.4. (Orthogonal tensor) [15] A tensor $\mathcal{Q} \in \mathbb{R}^{n \times n \times n_3}$ is orthogonal if it satisfies $\mathcal{Q}^* * \mathcal{Q} = \mathcal{Q} * \mathcal{Q}^* = \mathcal{I}$.

Definition 2.5. (F-diagonal Tensor) [15] A tensor is called f-diagonal if each of its frontal slices is a diagonal matrix.

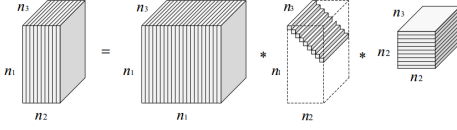


Fig. 2: An illustration of the t-SVD of an $n_1 \times n_2 \times n_3$ tensor [10].

Theorem 2.2. (T-SVD) Let $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$. Then it can be factorized as

$$\mathcal{A} = \mathcal{U} * \mathcal{S} * \mathcal{V}^*, \quad (16)$$

where $\mathcal{U} \in \mathbb{R}^{n_1 \times n_1 \times n_3}$, $\mathcal{V} \in \mathbb{R}^{n_2 \times n_2 \times n_3}$ are orthogonal, and $\mathcal{S} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ is an f -diagonal tensor.

Proof. The proof is by construction. Recall that (10) holds and $\bar{\mathcal{A}}^{(i)}$'s satisfy the property (11). Then we construct the SVD of each $\bar{\mathcal{A}}^{(i)}$ in the following way. For $i = 1, \dots, \lceil \frac{n_3+1}{2} \rceil$, let $\bar{\mathcal{A}}^{(i)} = \bar{\mathcal{U}}^{(i)} \bar{\mathcal{S}}^{(i)} (\bar{\mathcal{V}}^{(i)})^*$ be the full SVD of $\bar{\mathcal{A}}^{(i)}$. Here the singular values in $\bar{\mathcal{S}}^{(i)}$ are real. For $i = \lceil \frac{n_3+1}{2} \rceil + 1, \dots, n_3$, let $\bar{\mathcal{U}}^{(i)} = \text{conj}(\bar{\mathcal{U}}^{(n_3-i+2)})$, $\bar{\mathcal{S}}^{(i)} = \bar{\mathcal{S}}^{(n_3-i+2)}$ and $\bar{\mathcal{V}}^{(i)} = \text{conj}(\bar{\mathcal{V}}^{(n_3-i+2)})$. Then, it is easy to verify that $\bar{\mathcal{A}}^{(i)} = \bar{\mathcal{U}}^{(i)} \bar{\mathcal{S}}^{(i)} (\bar{\mathcal{V}}^{(i)})^*$ gives the full SVD of $\bar{\mathcal{A}}^{(i)}$ for $i = \lceil \frac{n_3+1}{2} \rceil + 1, \dots, n_3$. Then,

$$\bar{\mathcal{A}} = \bar{\mathcal{U}} \bar{\mathcal{S}} \bar{\mathcal{V}}^*. \quad (17)$$

By the construction of $\bar{\mathcal{U}}$, $\bar{\mathcal{S}}$ and $\bar{\mathcal{V}}$, and Lemma 2.1, we have that $(\mathbf{F}_{n_3}^{-1} \otimes \mathbf{I}_{n_1}) \cdot \bar{\mathcal{U}} \cdot (\mathbf{F}_{n_3} \otimes \mathbf{I}_{n_1})$, $(\mathbf{F}_{n_3}^{-1} \otimes \mathbf{I}_{n_1}) \cdot \bar{\mathcal{S}} \cdot (\mathbf{F}_{n_3} \otimes \mathbf{I}_{n_2})$ and $(\mathbf{F}_{n_3}^{-1} \otimes \mathbf{I}_{n_2}) \cdot \bar{\mathcal{V}} \cdot (\mathbf{F}_{n_3} \otimes \mathbf{I}_{n_2})$ are real block circulant matrices. Then we can obtain an expression for $\text{bcirc}(\mathcal{A})$ by applying the appropriate matrix $(\mathbf{F}_{n_3}^{-1} \otimes \mathbf{I}_{n_1})$ to the left and the appropriate matrix $(\mathbf{F}_{n_3} \otimes \mathbf{I}_{n_2})$ to the right of each of the matrices in (17), and folding up the result. This gives a decomposition of the form $\mathcal{U} * \mathcal{S} * \mathcal{V}^*$, where \mathcal{U} , \mathcal{S} and \mathcal{V} are real. \square

Theorem 2.2 shows that any 3 way tensor can be factorized into 3 components, including 2 orthogonal tensors and an f -diagonal tensor. See Figure 2 for an intuitive illustration of the t-SVD factorization. T-SVD reduces to the matrix SVD when $n_3 = 1$. We would like to emphasize that the result of Theorem 2.2 was given first in [15] and later in some other related works [10], [22]. But their proof and the way for computing \mathcal{U} and \mathcal{V} are not rigorous. The issue is that their method cannot guarantee that \mathcal{U} and \mathcal{V} are real tensors. They construct each frontal slice $\bar{\mathcal{U}}^{(i)}$ (or $\bar{\mathcal{V}}^{(i)}$) of $\bar{\mathcal{U}}$ (or $\bar{\mathcal{V}}$ resp.) from the SVD of $\bar{\mathcal{A}}^{(i)}$ independently for all $i = 1, \dots, n_3$. However, the matrix SVD is not unique. Thus, $\bar{\mathcal{U}}^{(i)}$'s and $\bar{\mathcal{V}}^{(i)}$'s may not satisfy property (11) even though $\bar{\mathcal{A}}^{(i)}$'s do. In this case, the obtained \mathcal{U} (or \mathcal{V}) from the inverse DFT of $\bar{\mathcal{U}}$ (or $\bar{\mathcal{V}}$ resp.) may not be real. Our proof above instead uses property (11) to construct \mathcal{U} and \mathcal{V} and thus avoids this issue. Our proof further leads to a more efficient way for computing t-SVD shown in Algorithm 2.

It is known that the singular values of a matrix have the decreasing order property. Let $\mathcal{A} = \mathcal{U} * \mathcal{S} * \mathcal{V}^*$ be the t-SVD of $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$. The entries on the diagonal of the first frontal slice $\mathcal{S}(:, :, 1)$ of \mathcal{S} have the same decreasing property, i.e.,

$$\mathcal{S}(1, 1, 1) \geq \mathcal{S}(2, 2, 1) \geq \dots \geq \mathcal{S}(n', n', 1) \geq 0, \quad (18)$$

Algorithm 2 T-SVD

Input: $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$.

Output: T-SVD components \mathcal{U} , \mathcal{S} and \mathcal{V} of \mathcal{A} .

1. Compute $\bar{\mathcal{A}} = \text{fft}(\mathcal{A}, [], 3)$.
2. Compute each frontal slice of $\bar{\mathcal{U}}$, $\bar{\mathcal{S}}$ and $\bar{\mathcal{V}}$ from $\bar{\mathcal{A}}$ by

for $i = 1, \dots, \lceil \frac{n_3+1}{2} \rceil$ **do**
 $[\bar{\mathcal{U}}^{(i)}, \bar{\mathcal{S}}^{(i)}, \bar{\mathcal{V}}^{(i)}] = \text{SVD}(\bar{\mathcal{A}}^{(i)});$
end for
for $i = \lceil \frac{n_3+1}{2} \rceil + 1, \dots, n_3$ **do**
 $\bar{\mathcal{U}}^{(i)} = \text{conj}(\bar{\mathcal{U}}^{(n_3-i+2)});$
 $\bar{\mathcal{S}}^{(i)} = \bar{\mathcal{S}}^{(n_3-i+2)};$
 $\bar{\mathcal{V}}^{(i)} = \text{conj}(\bar{\mathcal{V}}^{(n_3-i+2)});$
end for
3. Compute $\mathcal{U} = \text{ifft}(\bar{\mathcal{U}}, [], 3)$, $\mathcal{S} = \text{ifft}(\bar{\mathcal{S}}, [], 3)$, and $\mathcal{V} = \text{ifft}(\bar{\mathcal{V}}, [], 3)$.

where $n' = \min(n_1, n_2)$. The above property holds since the inverse DFT gives

$$\mathcal{S}(i, i, 1) = \frac{1}{n_3} \sum_{j=1}^{n_3} \bar{\mathcal{S}}(i, i, j), \quad (19)$$

and the entries on the diagonal of $\bar{\mathcal{S}}(:, :, j)$ are the singular values of $\bar{\mathcal{A}}(:, :, j)$. As will be seen in Section 3, the tensor nuclear norm depends only on the first frontal slice $\mathcal{S}(:, :, 1)$. Thus, we call the entries on the diagonal of $\mathcal{S}(:, :, 1)$ as the singular values of \mathcal{A} .

Definition 2.6. (Tensor tubal rank) [14], [34] For $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, the tensor tubal rank, denoted as $\text{rank}_t(\mathcal{A})$, is defined as the number of nonzero singular tubes of \mathcal{S} , where \mathcal{S} is from the t-SVD of $\mathcal{A} = \mathcal{U} * \mathcal{S} * \mathcal{V}^*$. We can write

$$\text{rank}_t(\mathcal{A}) = \#\{i, \mathcal{S}(i, i, :) \neq 0\}.$$

By using property (19), the tensor tubal rank is determined by the first frontal slice $\mathcal{S}(:, :, 1)$ of \mathcal{S} , i.e.,

$$\text{rank}_t(\mathcal{A}) = \#\{i, \mathcal{S}(i, i, 1) \neq 0\}.$$

Hence, the tensor tubal rank is equivalent to the number of nonzero singular values of \mathcal{A} . This property is the same as the matrix case. Define $\mathcal{A}_k = \sum_{i=1}^k \mathcal{U}(:, i, :) * \mathcal{S}(i, i, :) * \mathcal{V}(:, i, :)^*$ for some $k < \min(n_1, n_2)$. Then $\mathcal{A}_k = \arg \min_{\text{rank}_t(\tilde{\mathcal{A}}) \leq k} \|\mathcal{A} - \tilde{\mathcal{A}}\|_F$, so \mathcal{A}_k is the best approximation of \mathcal{A} with the tubal rank at most k . It is known that the real color images can be well approximated by low-rank matrices on the three channels independently. If we treat a color image as a three way tensor with each channel corresponding to a frontal slice, then it can be well approximated by a tensor of low tubal rank. A similar observation was found in [10] with the application to facial recognition. Figure 3 gives an example to show that a color image can be well approximated by a low tubal rank tensor since most of the singular values of the corresponding tensor are relatively small.

In Section 3, we will define a new tensor nuclear norm which is the convex surrogate of the tensor average rank defined as follows. This rank is closely related to the tensor tubal rank.

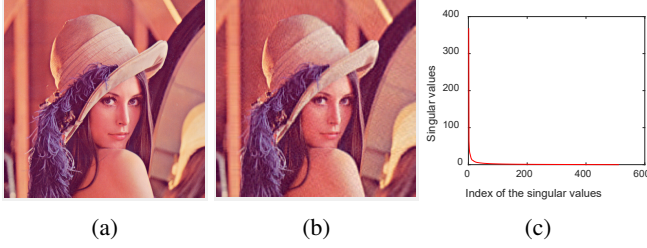


Fig. 3: Color images can be approximated by low tubal rank tensors. (a) A color image can be modeled as a tensor $\mathcal{M} \in \mathbb{R}^{512 \times 512 \times 3}$; (b) approximation by a tensor with tubal rank $r = 50$; (c) plot of the singular values of \mathcal{M} .

Definition 2.7. (Tensor average rank) For $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, the tensor average rank, denoted as $\text{rank}_a(\mathcal{A})$, is defined as

$$\text{rank}_a(\mathcal{A}) = \frac{1}{n_3} \text{rank}(\text{bcirc}(\mathcal{A})). \quad (20)$$

The above definition has a factor $\frac{1}{n_3}$. Note that this factor is crucial in this work as it guarantees that the convex envelope of the tensor average rank within a certain set is the tensor nuclear norm defined in Section 3. The underlying reason for this factor is the t-product definition. Each element of \mathcal{A} is repeated n_3 times in the block circulant matrix $\text{bcirc}(\mathcal{A})$ used in the t-product. Intuitively, this factor alleviates such an entries expansion issue.

There are some connections between different tensor ranks and these properties imply that the low tubal rank or low average rank assumptions are reasonable for their applications in real visual data. First, $\text{rank}_a(\mathcal{A}) \leq \text{rank}_t(\mathcal{A})$. Indeed,

$$\text{rank}_a(\mathcal{A}) = \frac{1}{n_3} \text{rank}(\bar{\mathcal{A}}) \leq \max_{i=1, \dots, n_3} \text{rank}(\bar{\mathcal{A}}^{(i)}) = \text{rank}_t(\mathcal{A}),$$

where the first equality uses (10). This implies that a low tubal rank tensor always has low average rank. Second, let $\text{rank}_{tc}(\mathcal{A}) = (\text{rank}(\mathcal{A}^{\{1\}}), \text{rank}(\mathcal{A}^{\{2\}}), \text{rank}(\mathcal{A}^{\{3\}}))$, where $\mathcal{A}^{\{i\}}$ is the mode- i matricization of \mathcal{A} , be the Tucker rank of \mathcal{A} . Then $\text{rank}_a(\mathcal{A}) \leq \text{rank}(\mathcal{A}^{\{1\}})$. This implies that a tensor with low Tucker rank has low average rank. The low Tucker rank assumption used in some applications, e.g., image completion [18], is applicable to the low average rank assumption. Third, if the CP rank of \mathcal{A} is r , then its tubal rank is at most r [33]. Let $\mathcal{A} = \sum_{i=1}^r \mathbf{a}_i^{(1)} \circ \mathbf{a}_i^{(2)} \circ \mathbf{a}_i^{(3)}$, where \circ denotes the outer product, be the CP decomposition of \mathcal{A} . Then $\bar{\mathcal{A}} = \sum_{i=1}^r \mathbf{a}_i^{(1)} \circ \mathbf{a}_i^{(2)} \circ \bar{\mathbf{a}}_i^{(3)}$, where $\bar{\mathbf{a}}_i^{(3)} = \text{fft}(\mathbf{a}_i^{(3)})$. So $\bar{\mathcal{A}}$ has the CP rank at most r , and each frontal slice of $\bar{\mathcal{A}}$ is the sum of r rank-1 matrices. Thus, the tubal rank of \mathcal{A} is at most r . In summary, we show that the low average rank assumption is weaker than the low Tucker rank and low CP rank assumptions.

3 TENSOR NUCLEAR NORM (TNN)

In this section, we propose a new tensor nuclear norm which is a convex surrogate of tensor average rank. Based on t-SVD, one may have many different ways to define the tensor nuclear norm intuitively. We give a new and rigorous way to deduce the

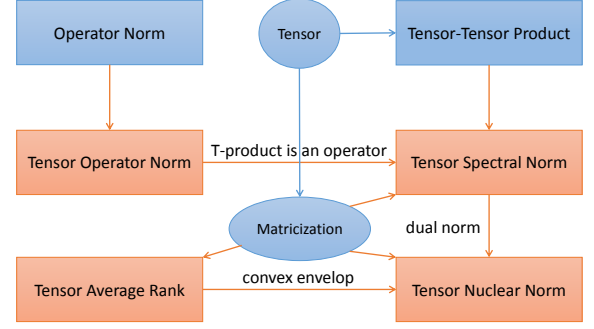


Fig. 4: An illustration of the way to define the tensor nuclear norm and the relationship with other tensor concepts. First, the tensor operator norm is a special case of the known operator norm performed on the tensors. The tensor spectral norm is induced by the tensor operator norm by treating the tensor-tensor product as an operator. Then the tensor nuclear norm is defined as the dual norm of the tensor spectral norm. We also define the tensor average rank and show that its convex envelope is the tensor nuclear norm within the unit ball of the tensor spectral norm. As detailed in Section 3, the tensor spectral norm, tensor nuclear norm and tensor average rank are also defined on the matricization of the tensor.

tensor nuclear norm from the t-product, such that the concepts and their properties are consistent with the matrix cases. This is important since it guarantees that the theoretical analysis of the tensor nuclear norm based tensor RPCA model in Section 4 can be done in a similar way to RPCA. Figure 4 summarizes the way for the new definitions and their relationships. It begins with the known operator norm [1] and t-product. First, the tensor spectral norm is induced by the tensor operator norm by treating the t-product as an operator. Then the tensor nuclear norm is defined as the dual norm of the tensor spectral norm. Finally, we show that the tensor nuclear norm is the convex envelope of the tensor average rank within the unit ball of the tensor spectral norm.

Let us first recall the concept of operator norm [1]. Let $(V, \|\cdot\|_V)$ and $(W, \|\cdot\|_W)$ be normed linear spaces and $L : V \rightarrow W$ be the bounded linear operator between them, respectively. The operator norm of L is defined as

$$\|L\| = \sup_{\|v\|_V \leq 1} \|L(v)\|_W. \quad (21)$$

Let $V = \mathbb{C}^{n_2}$, $W = \mathbb{C}^{n_1}$ and $L(v) = \mathcal{A}v$, $v \in V$, where $\mathcal{A} \in \mathbb{C}^{n_1 \times n_2}$. Based on different choices of $\|\cdot\|_V$ and $\|\cdot\|_W$, many matrix norms can be induced by the operator norm in (21). For example, if $\|\cdot\|_V$ and $\|\cdot\|_W$ are $\|\cdot\|_F$, then the operator norm (21) reduces to the matrix spectral norm.

Now, consider the normed linear spaces $(V, \|\cdot\|_F)$ and $(W, \|\cdot\|_F)$, where $V = \mathbb{R}^{n_2 \times 1 \times n_3}$, $W = \mathbb{R}^{n_1 \times 1 \times n_3}$, and $\mathcal{L} : V \rightarrow W$ is a bounded linear operator. In this case, (21) reduces to the tensor operator norm

$$\|\mathcal{L}\| = \sup_{\|\mathcal{V}\|_F \leq 1} \|\mathcal{L}(\mathcal{V})\|_F. \quad (22)$$

As a special case, if $\mathcal{L}(\mathcal{V}) = \mathcal{A} * \mathcal{V}$, where $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ and $\mathcal{V} \in V$, then the tensor operator norm (22) gives the tensor

spectral norm, denoted as $\|\mathcal{A}\|$,

$$\begin{aligned} \|\mathcal{A}\| &:= \sup_{\|\mathcal{V}\|_F \leq 1} \|\mathcal{A} * \mathcal{V}\|_F \\ &= \sup_{\|\mathcal{V}\|_F \leq 1} \|\text{bcirc}(\mathcal{A}) \cdot \text{unfold}(\mathcal{V})\|_F \end{aligned} \quad (23)$$

$$= \|\text{bcirc}(\mathcal{A})\|, \quad (24)$$

where (23) uses (14), and (24) uses the definition of matrix spectral norm.

Definition 3.1. (Tensor spectral norm) The tensor spectral norm of $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ is defined as $\|\mathcal{A}\| := \|\text{bcirc}(\mathcal{A})\|$.

By (7) and (10), we have

$$\|\mathcal{A}\| = \|\text{bcirc}(\mathcal{A})\| = \|\bar{\mathcal{A}}\|. \quad (25)$$

This property is frequently used in this work. It is known that the matrix nuclear norm is the dual norm of the matrix spectral norm. Thus, we define the tensor nuclear norm, denoted as $\|\mathcal{A}\|_*$, as the dual norm of the tensor spectral norm. For any $\mathcal{B} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ and $\bar{\mathcal{B}} \in \mathbb{C}^{n_1 n_3 \times n_2 n_3}$, we have

$$\|\mathcal{A}\|_* := \sup_{\|\mathcal{B}\| \leq 1} \langle \mathcal{A}, \mathcal{B} \rangle \quad (26)$$

$$= \sup_{\|\bar{\mathcal{B}}\| \leq 1} \frac{1}{n_3} \langle \bar{\mathcal{A}}, \bar{\mathcal{B}} \rangle \quad (27)$$

$$\leq \frac{1}{n_3} \sup_{\|\bar{\mathcal{B}}\| \leq 1} |\langle \bar{\mathcal{A}}, \bar{\mathcal{B}} \rangle| \quad (28)$$

$$= \frac{1}{n_3} \|\bar{\mathcal{A}}\|_*, \quad (29)$$

$$= \frac{1}{n_3} \|\text{bcirc}(\mathcal{A})\|_*, \quad (30)$$

where (27) is from (13), (28) is due to the fact that $\bar{\mathcal{B}}$ is a block diagonal matrix in $\mathbb{C}^{n_1 n_3 \times n_2 n_3}$ while $\bar{\mathcal{B}}$ is an arbitrary matrix in $\mathbb{C}^{n_1 n_3 \times n_2 n_3}$, (29) uses the fact that the matrix nuclear norm is the dual norm of the matrix spectral norm, and (30) uses (10) and (7). Now we show that there exists $\mathcal{B} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ such that the equality (28) holds and thus $\|\mathcal{A}\|_* = \frac{1}{n_3} \|\text{bcirc}(\mathcal{A})\|_*$. Let $\mathcal{A} = \mathcal{U} * \mathcal{S} * \mathcal{V}^*$ be the t-SVD of \mathcal{A} and $\mathcal{B} = \mathcal{U} * \mathcal{V}^*$. We have

$$\begin{aligned} \langle \mathcal{A}, \mathcal{B} \rangle &= \langle \mathcal{U} * \mathcal{S} * \mathcal{V}^*, \mathcal{U} * \mathcal{V}^* \rangle \\ &= \frac{1}{n_3} \langle \overline{\mathcal{U} * \mathcal{S} * \mathcal{V}^*}, \overline{\mathcal{U} * \mathcal{V}^*} \rangle \\ &= \frac{1}{n_3} \langle \bar{\mathcal{U}} \bar{\mathcal{S}} \bar{\mathcal{V}}^*, \bar{\mathcal{U}} \bar{\mathcal{V}}^* \rangle = \frac{1}{n_3} \text{Tr}(\bar{\mathcal{S}}) \\ &= \frac{1}{n_3} \|\bar{\mathcal{A}}\|_* = \frac{1}{n_3} \|\text{bcirc}(\mathcal{A})\|_*. \end{aligned} \quad (31)$$

Combining (26)-(30) and (31)-(32) leads to $\|\mathcal{A}\|_* = \frac{1}{n_3} \|\text{bcirc}(\mathcal{A})\|_*$. On the other hand, by (31)-(32), we have

$$\begin{aligned} \|\mathcal{A}\|_* &= \langle \mathcal{U} * \mathcal{S} * \mathcal{V}^*, \mathcal{U} * \mathcal{V}^* \rangle \\ &= \langle \mathcal{U}^* * \mathcal{U} * \mathcal{S}, \mathcal{V}^* * \mathcal{V} \rangle \\ &= \langle \mathcal{S}, \mathcal{I} \rangle = \sum_{i=1}^r \mathcal{S}(i, i, 1), \end{aligned} \quad (32)$$

where $r = \text{rank}_t(\mathcal{A})$ is the tubal rank. Thus, we have the following definition of tensor nuclear norm.

Definition 3.2. (Tensor nuclear norm) Let $\mathcal{A} = \mathcal{U} * \mathcal{S} * \mathcal{V}^*$ be the t-SVD of $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$. The tensor nuclear norm of \mathcal{A} is defined as

$$\|\mathcal{A}\|_* := \langle \mathcal{S}, \mathcal{I} \rangle = \sum_{i=1}^r \mathcal{S}(i, i, 1),$$

where $r = \text{rank}_t(\mathcal{A})$.

From (33), it can be seen that only the information in the first frontal slice of \mathcal{S} is used when defining the tensor nuclear norm. Note that this is the first work which directly uses the singular values $\mathcal{S}(:, :, 1)$ of a tensor to define the tensor nuclear norm. Such a definition makes it consistent with the matrix nuclear norm. The above TNN definition is also different from existing works [19], [34], [27].

It is known that the matrix nuclear norm $\|\mathcal{A}\|_*$ is the convex envelope of the matrix rank $\text{rank}(\mathcal{A})$ within the set $\{\mathcal{A} | \|\mathcal{A}\| \leq 1\}$ [9]. Now we show that the tensor average rank and tensor nuclear norm have the same relationship.

Theorem 3.1. On the set $\{\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3} | \|\mathcal{A}\| \leq 1\}$, the convex envelope of the tensor average rank $\text{rank}_a(\mathcal{A})$ is the tensor nuclear norm $\|\mathcal{A}\|_*$.

We would like to emphasize that the proposed tensor spectral norm, tensor nuclear norm and tensor ranks are not arbitrarily defined. They are rigorously induced by the t-product and t-SVD. These concepts and their relationships are consistent with the matrix cases. This is important for the proofs, analysis and computation in optimization. Table 1 summarizes the parallel concepts in sparse vector, low-rank matrix and low-rank tensor. With these elements in place, the existing proofs of low-rank matrix recovery provide a template for the more general case of low-rank tensor recovery.

Also, from the above discussions, we have the property

$$\|\mathcal{A}\|_* = \frac{1}{n_3} \|\text{bcirc}(\mathcal{A})\|_* = \frac{1}{n_3} \|\bar{\mathcal{A}}\|_*. \quad (34)$$

It is interesting to understand the tensor nuclear norm from the perspectives of $\text{bcirc}(\mathcal{A})$ and $\bar{\mathcal{A}}$. The block circulant matrix can be regarded as a new way of matricization of \mathcal{A} in the original domain. The frontal slices of \mathcal{A} are arranged in a circulant way, which is expected to preserve more spacial relationships across frontal slices, compared with previous matricizations along a single dimension. Also, the block diagonal matrix $\bar{\mathcal{A}}$ can be regarded as a matricization of \mathcal{A} in the Fourier domain. Its blocks on the diagonal are the frontal slices of $\bar{\mathcal{A}}$, which contains the information across frontal slices of \mathcal{A} due to the DFT on \mathcal{A} along the third dimension. So $\text{bcirc}(\mathcal{A})$ and $\bar{\mathcal{A}}$ play a similar role to matricizations of \mathcal{A} in different domains. Both of them capture the spacial information within and across frontal slices of \mathcal{A} . This intuitively supports our tensor nuclear norm definition.

Let $\mathcal{A} = \mathcal{U} \mathcal{S} \mathcal{V}^*$ be the skinny SVD of \mathcal{A} . It is known that any subgradient of the nuclear norm at \mathcal{A} is of the form $\mathcal{U} \mathcal{V}^* + \mathcal{W}$, where $\mathcal{U}^* \mathcal{W} = \mathbf{0}$, $\mathcal{W} \mathcal{V} = \mathbf{0}$ and $\|\mathcal{W}\| \leq 1$ [32].

TABLE 1: Parallelism of sparse vector, low-rank matrix and low-rank tensor.

	Sparse vector	Low-rank matrix	Low-rank tensor (this work)
Degeneracy of	1-D signal $\mathbf{x} \in \mathbb{R}^n$	2-D correlated signals $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$	3-D correlated signals $\mathcal{X} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$
Parsimony concept	cardinality	rank	tensor average rank ¹
Measure	ℓ_0 -norm $\ \mathbf{x}\ _0$	$\text{rank}(\mathbf{X})$	$\text{rank}_a(\mathcal{X})$
Convex surrogate	ℓ_1 -norm $\ \mathbf{x}\ _1$	nuclear norm $\ \mathbf{X}\ _*$	tensor nuclear norm $\ \mathcal{X}\ _*$
Dual norm	ℓ_∞ -norm $\ \mathbf{x}\ _\infty$	spectral norm $\ \mathbf{X}\ $	tensor spectral norm $\ \mathcal{X}\ $

¹Strictly speaking, the tensor tubal rank, which bounds the tensor average rank, is also the parsimony concept of the low-rank tensor.

Similarly, for $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ with tubal rank r , we also have the skinny t-SVD, i.e., $\mathcal{A} = \mathcal{U} * \mathcal{S} * \mathcal{V}^*$, where $\mathcal{U} \in \mathbb{R}^{n_1 \times r \times n_3}$, $\mathcal{S} \in \mathbb{R}^{r \times r \times n_3}$, and $\mathcal{V} \in \mathbb{R}^{n_2 \times r \times n_3}$, in which $\mathcal{U}^* * \mathcal{U} = \mathcal{I}$ and $\mathcal{V}^* * \mathcal{V} = \mathcal{I}$. The skinny t-SVD will be used throughout this paper. With skinny t-SVD, we introduce the subgradient of the tensor nuclear norm, which plays an important role in the proofs.

Theorem 3.2. (Subgradient of tensor nuclear norm) Let $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ with $\text{rank}_t(\mathcal{A}) = r$ and its skinny t-SVD be $\mathcal{A} = \mathcal{U} * \mathcal{S} * \mathcal{V}^*$. The subdifferential (the set of subgradients) of $\|\mathcal{A}\|_*$ is $\partial\|\mathcal{A}\|_* = \{\mathcal{U} * \mathcal{V}^* + \mathcal{W} | \mathcal{U}^* * \mathcal{W} = \mathbf{0}, \mathcal{W} * \mathcal{V} = \mathbf{0}, \|\mathcal{W}\| \leq 1\}$.

4 EXACT RECOVERY GUARANTEE OF TRPCA

With TNN defined above, we now consider the exact recovery guarantee of TRPCA in (5). The problem we study here is to recover a low tubal rank tensor \mathcal{L}_0 from highly corrupted measurements $\mathcal{X} = \mathcal{L}_0 + \mathcal{S}_0$. In this section, we show that under certain assumptions, the low tubal rank part \mathcal{L}_0 and sparse part \mathcal{S}_0 can be exactly recovered by solving convex program (5). We will also give the optimization detail for solving (5).

4.1 Tensor Incoherence Conditions

Recovering the low-rank and sparse components from their sum suffers from an identifiability issue. For example, the tensor \mathcal{X} , with $x_{ijk} = 1$ when $i = j = k = 1$ and zeros everywhere else, is both low-rank and sparse. One is not able to identify the low-rank component and the sparse component in this case. To avoid such pathological situations, we need to assume that the low-rank component \mathcal{L}_0 is not sparse. To this end, we assume \mathcal{L}_0 to satisfy some incoherence conditions. We denote $\hat{\mathbf{e}}_i$ as the tensor column basis, which is a tensor of size $n_1 \times 1 \times n_3$ with its $(i, 1, 1)$ -th entry equaling 1 and the rest equaling 0 [33]. We also define the tensor tube basis $\hat{\mathbf{e}}_k$, which is a tensor of size $1 \times 1 \times n_3$ with its $(1, 1, k)$ -th entry equaling 1 and the rest equaling 0.

Definition 4.1. (Tensor Incoherence Conditions) For $\mathcal{L}_0 \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, assume that $\text{rank}_t(\mathcal{L}_0) = r$ and it has the skinny t-SVD $\mathcal{L}_0 = \mathcal{U} * \mathcal{S} * \mathcal{V}^*$, where $\mathcal{U} \in \mathbb{R}^{n_1 \times r \times n_3}$ and $\mathcal{V} \in \mathbb{R}^{n_2 \times r \times n_3}$ satisfy $\mathcal{U}^* * \mathcal{U} = \mathcal{I}$ and $\mathcal{V}^* * \mathcal{V} = \mathcal{I}$, and $\mathcal{S} \in \mathbb{R}^{r \times r \times n_3}$ is an f -diagonal tensor. Then \mathcal{L}_0 is said to satisfy the tensor incoherence conditions with parameter μ if

$$\max_{i=1, \dots, n_1} \|\mathcal{U}^* * \hat{\mathbf{e}}_i\|_F \leq \sqrt{\frac{\mu r}{n_1 n_3}}, \quad (35)$$

$$\max_{j=1, \dots, n_2} \|\mathcal{V}^* * \hat{\mathbf{e}}_j\|_F \leq \sqrt{\frac{\mu r}{n_2 n_3}}, \quad (36)$$

$$\|\mathcal{U} * \mathcal{V}^*\|_\infty \leq \sqrt{\frac{\mu r}{n_1 n_2 n_3^2}}. \quad (37)$$

The exact recovery guarantee of RPCA [3] also requires some incoherence conditions. Due to property (12), conditions (48)-(49) have equivalent matrix forms in the Fourier domain, and they are intuitively similar to the matrix incoherence conditions (1.2) in [3]. But the joint incoherence condition (50) is more different from the matrix case (1.3) in [3], since it does not have an equivalent matrix form in the Fourier domain. As observed in [4], the joint incoherence condition is not necessary for low-rank matrix completion. However, for RPCA, it is unavoidable for polynomial-time algorithms. In our proofs, the joint incoherence (50) condition is necessary. Another identifiability issue arises if the sparse tensor \mathcal{S}_0 has low tubal rank. This can be avoided by assuming that the support of \mathcal{S}_0 is uniformly distributed.

4.2 Main Results

Now we show that the convex program (5) is able to perfectly recover the low-rank and sparse components. We define $n_{(1)} = \max(n_1, n_2)$ and $n_{(2)} = \min(n_1, n_2)$.

Theorem 4.1. Suppose that $\mathcal{L}_0 \in \mathbb{R}^{n \times n \times n_3}$ obeys (48)-(50). Fix any $n \times n \times n_3$ tensor \mathcal{M} of signs. Suppose that the support set Ω of \mathcal{S}_0 is uniformly distributed among all sets of cardinality m , and that $\text{sgn}([\mathcal{S}_0]_{ijk}) = [\mathcal{M}]_{ijk}$ for all $(i, j, k) \in \Omega$. Then, there exist universal constants $c_1, c_2 > 0$ such that with probability at least $1 - c_1(n n_3)^{-c_2}$ (over the choice of support of \mathcal{S}_0), $(\mathcal{L}_0, \mathcal{S}_0)$ is the unique minimizer to (5) with $\lambda = 1/\sqrt{n n_3}$, provided that

$$\text{rank}_t(\mathcal{L}_0) \leq \frac{\rho_r n n_3}{\mu(\log(n n_3))^2} \text{ and } m \leq \rho_s n^2 n_3, \quad (38)$$

where ρ_r and ρ_s are positive constants. If $\mathcal{L}_0 \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ has rectangular frontal slices, TRPCA with $\lambda = 1/\sqrt{n_{(1)} n_3}$ succeeds with probability at least $1 - c_1(n_{(1)} n_3)^{-c_2}$, provided that $\text{rank}_t(\mathcal{L}_0) \leq \frac{\rho_r n_{(2)} n_3}{\mu(\log(n_{(1)} n_3))^2}$ and $m \leq \rho_s n_1 n_2 n_3$.

The above result shows that for incoherent \mathcal{L}_0 , the perfect recovery is guaranteed with high probability for $\text{rank}_t(\mathcal{L}_0)$ on the order of $n n_3 / (\mu(\log n n_3)^2)$ and a number of nonzero entries in \mathcal{S}_0 on the order of $n^2 n_3$. For \mathcal{S}_0 , we make only one assumption on the random location distribution, but no assumption about the magnitudes or signs of the nonzero entries. Also TRPCA is parameter free. The mathematical analysis implies that the parameter $\lambda = 1/\sqrt{n n_3}$ leads to the correct recovery. Moreover, since the t-product of 3-way tensors reduces to the standard matrix-matrix product when $n_3 = 1$, the tensor nuclear norm reduces to the matrix nuclear norm. Thus, RPCA is a special case of TRPCA and the guarantee of RPCA in Theorem 1.1 in [3] is a special case of our Theorem 4.1. Both our model and theoretical

guarantee are consistent with RPCA. Compared with SNN [13], our tensor extension of RPCA is much more simple and elegant.

The detailed proof of Theorem 4.1 can be found in the supplementary material. It is interesting to understand our proof from the perspective of the following equivalent formulation

$$\min_{\mathcal{L}, \mathcal{E}} \frac{1}{n_3} (\|\bar{\mathcal{L}}\|_* + \lambda \|\text{bcirc}(\mathcal{E})\|_1), \text{ s.t. } \mathcal{X} = \mathcal{L} + \mathcal{E}, \quad (39)$$

where (34) is used. Program (39) is a mixed model since the low-rank regularization is performed on the Fourier domain while the sparse regularization is performed on the original domain. Our proof of Theorem 4.1 is also conducted based on the interaction between both domains. By interpreting the tensor nuclear norm of \mathcal{L} as the matrix nuclear norm of $\bar{\mathcal{L}}$ (with a factor $\frac{1}{n_3}$) in the Fourier domain, we are then able to use some existing properties of the matrix nuclear norm in the proofs. The analysis for the sparse term is kept on the original domain since the ℓ_1 -norm has no equivalent form in the Fourier domain. Though both two terms of the objective function of (39) are given on two matrices ($\bar{\mathcal{L}}$ and $\text{bcirc}(\mathcal{E})$), the analysis for model (39) is very different from that of matrix RPCA. The matrices $\bar{\mathcal{L}}$ and $\text{bcirc}(\mathcal{E})$ can be regarded as two matricizations of the tensor objects \mathcal{L} and \mathcal{E} , respectively. Their structures are more complicated than those in matrix RPCA, and thus make the proofs different from [3]. For example, our proofs require proving several bounds of norms on random tensors. These results and proofs, which have to use the properties of block circulant matrices and the Fourier transformation, are completely new. Some proofs are challenging due to the dependent structure of $\text{bcirc}(\mathcal{E})$ for \mathcal{E} with an independent elements assumption. Also, TRPCA is of a different nature from the tensor completion problem [33]. The proof of the exact recovery of TRPCA is more challenging since the ℓ_1 -norm (and its dual norm ℓ_∞ -norm used in (50)) has no equivalent formulation in the Fourier domain.

It is worth mentioning that this work focuses on the analysis for 3-way tensors. But it is not difficult to generalize our model in (5) and results in Theorem 4.1 to the case of order- p ($p \geq 3$) tensors, by using the t-SVD for order- p tensors in [22].

When considering the application of TRPCA, the way for constructing a 3-way tensor from data is important. The reason is that the t-product is orientation dependent, and so is the tensor nuclear norm. Thus, the value of TNN may be different if the tensor is rotated. For example, a 3-channel color image can be formatted as 3 different sizes of tensors. Therefore, when using TRPCA which is based on TNN, one has to format the data into tensors in a proper way by leveraging some priori knowledge, e.g., the low tubal rank property of the constructed tensor.

4.3 Tensor Singular Value Thresholding

Problem (5) can be solved by the standard Alternating Direction Method of Multiplier (ADMM) [20]. A key step is to compute the proximal operator of TNN

$$\min_{\mathcal{X} \in \mathbb{R}^{n_1 \times n_2 \times n_3}} \tau \|\mathcal{X}\|_* + \frac{1}{2} \|\mathcal{X} - \mathcal{Y}\|_F^2. \quad (40)$$

We show that it also has a closed-form solution as the proximal operator of the matrix nuclear norm. Let $\mathcal{Y} = \mathcal{U} * \mathcal{S} * \mathcal{V}^*$ be the

Algorithm 3 Tensor Singular Value Thresholding (t-SVT)

Input: $\mathcal{Y} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, $\tau > 0$.

Output: $\mathcal{D}_\tau(\mathcal{Y})$ as defined in (41).

1. Compute $\bar{\mathcal{Y}} = \text{fft}(\mathcal{Y}, [], 3)$.
 2. Perform matrix SVT on each frontal slice of $\bar{\mathcal{Y}}$ by

for $i = 1, \dots, \lceil \frac{n_3+1}{2} \rceil$ **do**
 $[\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{SVD}(\bar{\mathcal{Y}}^{(i)})$;
 $\bar{\mathcal{W}}^{(i)} = \mathbf{U} \cdot (\mathbf{S} - \tau)_+ \cdot \mathbf{V}^*$;
end for
for $i = \lceil \frac{n_3+1}{2} \rceil + 1, \dots, n_3$ **do**
 $\bar{\mathcal{W}}^{(i)} = \text{conj}(\bar{\mathcal{W}}^{(n_3-i+2)})$;
end for
 3. Compute $\mathcal{D}_\tau(\mathcal{Y}) = \text{ifft}(\bar{\mathcal{W}}, [], 3)$.
-

tensor SVD of $\mathcal{Y} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$. For each $\tau > 0$, we define the tensor Singular Value Thresholding (t-SVT) operator as follows

$$\mathcal{D}_\tau(\mathcal{Y}) = \mathcal{U} * \mathcal{S}_\tau * \mathcal{V}^*, \quad (41)$$

where

$$\mathcal{S}_\tau = \text{ifft}((\bar{\mathcal{S}} - \tau)_+, [], 3). \quad (42)$$

Note that $\bar{\mathcal{S}}$ is a real tensor. Above t_+ denotes the positive part of t , i.e., $t_+ = \max(t, 0)$. That is, this operator simply applies a soft-thresholding rule to the singular values $\bar{\mathcal{S}}$ (not \mathcal{S}) of the frontal slices of $\bar{\mathcal{Y}}$, effectively shrinking these towards zero. The t-SVT operator is the proximity operator associated with TNN.

Theorem 4.2. For any $\tau > 0$ and $\mathcal{Y} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, the tensor singular value thresholding operator (41) obeys

$$\mathcal{D}_\tau(\mathcal{Y}) = \arg \min_{\mathcal{X} \in \mathbb{R}^{n_1 \times n_2 \times n_3}} \tau \|\mathcal{X}\|_* + \frac{1}{2} \|\mathcal{X} - \mathcal{Y}\|_F^2. \quad (43)$$

Proof. The required solution to (43) is a real tensor and thus we first show that $\mathcal{D}_\tau(\mathcal{Y})$ in (41) is real. Let $\mathcal{Y} = \mathcal{U} * \mathcal{S} * \mathcal{V}^*$ be the tensor SVD of \mathcal{Y} . We know that the frontal slices of $\bar{\mathcal{S}}$ satisfy the property (11) and so do the frontal slices of $(\bar{\mathcal{S}} - \tau)_+$. By Lemma 2.1, \mathcal{S}_τ in (42) is real. Thus, $\mathcal{D}_\tau(\mathcal{Y})$ in (41) is real. Secondly, by using properties (34) and (12), problem (43) is equivalent to

$$\begin{aligned} & \arg \min_{\mathcal{X}} \frac{1}{n_3} (\tau \|\bar{\mathcal{X}}\|_* + \frac{1}{2} \|\bar{\mathcal{X}} - \bar{\mathcal{Y}}\|_F^2) \\ &= \arg \min_{\mathcal{X}} \frac{1}{n_3} \sum_{i=1}^{n_3} (\tau \|\bar{\mathcal{X}}^{(i)}\|_* + \frac{1}{2} \|\bar{\mathcal{X}}^{(i)} - \bar{\mathcal{Y}}^{(i)}\|_F^2). \end{aligned} \quad (44)$$

By Theorem 2.1 in [2], we know that the i -th frontal slice of $\bar{\mathcal{D}}_\tau(\mathcal{Y})$ solves the i -th subproblem of (44). Hence, $\mathcal{D}_\tau(\mathcal{Y})$ solves problem (43). \square

Theorem 4.2 gives the closed-form of the t-SVT operator $\mathcal{D}_\tau(\mathcal{Y})$, which is a natural extension of the matrix SVT [2]. Note that $\mathcal{D}_\tau(\mathcal{Y})$ is real when \mathcal{Y} is real. By using property (11), Algorithm 3 gives an efficient way for computing $\mathcal{D}_\tau(\mathcal{Y})$.

With t-SVT, we now give the details of ADMM to solve (5). The augmented Lagrangian function of (5) is

$$\begin{aligned} L(\mathcal{L}, \mathcal{E}, \mathcal{Y}, \mu) &= \|\mathcal{L}\|_* + \lambda \|\mathcal{E}\|_1 + \langle \mathcal{Y}, \mathcal{L} + \mathcal{E} - \mathcal{X} \rangle \\ &\quad + \frac{\mu}{2} \|\mathcal{L} + \mathcal{E} - \mathcal{X}\|_F^2. \end{aligned}$$

Algorithm 4 Solve (5) by ADMM

Initialize: $\mathcal{L}_0 = \mathcal{S}_0 = \mathcal{Y}_0 = 0$, $\rho = 1.1$, $\mu_0 = 1e-3$, $\mu_{\max} = 1e10$, $\epsilon = 1e-8$.

while not converged **do**

1. Update \mathcal{L}_{k+1} by

$$\mathcal{L}_{k+1} = \underset{\mathcal{L}}{\operatorname{argmin}} \left\| \mathcal{L} \right\|_* + \frac{\mu_k}{2} \left\| \mathcal{L} + \mathcal{E}_k - \mathcal{X} + \frac{\mathcal{Y}_k}{\mu_k} \right\|_F^2;$$

2. Update \mathcal{E}_{k+1} by

$$\mathcal{E}_{k+1} = \underset{\mathcal{E}}{\operatorname{argmin}} \lambda \|\mathcal{E}\|_1 + \frac{\mu_k}{2} \left\| \mathcal{L}_{k+1} + \mathcal{E} - \mathcal{X} + \frac{\mathcal{Y}_k}{\mu_k} \right\|_F^2;$$

3. $\mathcal{Y}_{k+1} = \mathcal{Y}_k + \mu_k(\mathcal{L}_{k+1} + \mathcal{E}_{k+1} - \mathcal{X})$;

4. Update μ_{k+1} by $\mu_{k+1} = \min(\rho\mu_k, \mu_{\max})$;

5. Check the convergence conditions

$$\begin{aligned} \|\mathcal{L}_{k+1} - \mathcal{L}_k\|_\infty &\leq \epsilon, \quad \|\mathcal{E}_{k+1} - \mathcal{E}_k\|_\infty \leq \epsilon, \\ \|\mathcal{L}_{k+1} + \mathcal{E}_{k+1} - \mathcal{X}\|_\infty &\leq \epsilon. \end{aligned}$$

end while

Then \mathcal{L} and \mathcal{E} can be updated by minimizing the augmented Lagrangian function L alternately. Both subproblems have closed-form solutions. See Algorithm 4 for the whole procedure. The main per-iteration cost lies in the update of \mathcal{L}_{k+1} , which requires computing FFT and $\lceil \frac{n_3+1}{2} \rceil$ SVDs of $n_1 \times n_2$ matrices. The per-iteration complexity is $O(n_1 n_2 n_3 \log n_3 + n_{(1)} n_{(2)}^2 n_3)$.

5 EXPERIMENTS

In this section, we conduct numerical experiments to verify our main results in Theorem 4.1. We first investigate the ability of the convex TRPCA model (5) to recover tensors with varying tubal rank and different levels of sparse noises. We then apply it for image recovery and background modeling. As suggested by Theorem 4.1, we set $\lambda = 1/\sqrt{n_{(1)}n_3}$ in all the experiments². But note that it is possible to further improve the performance by tuning λ more carefully. The suggested value in theory provides a good guide in practice. All the simulations are conducted on a PC with an Intel Xeon E3-1270 3.60GHz CPU and 64GB memory.

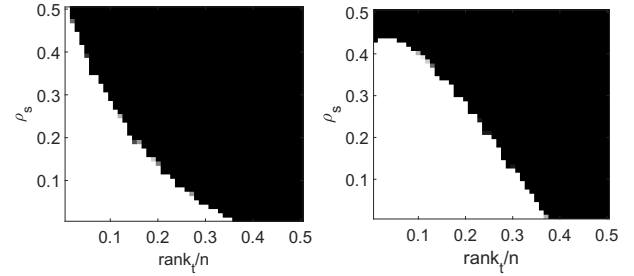
5.1 Exact Recovery from Varying Fractions of Error

We first verify the correct recovery guarantee of Theorem 4.1 on randomly generated problems. We simply consider the tensors of size $n \times n \times n$, with varying dimension $n = 100$ and 200 . We generate a tensor with tubal rank r as a product $\mathcal{L}_0 = \mathcal{P} * \mathcal{Q}^*$, where \mathcal{P} and \mathcal{Q} are $n \times r \times n$ tensors with entries independently sampled from $\mathcal{N}(0, 1/n)$ distribution. The support set Ω (with size m) of \mathcal{E}_0 is chosen uniformly at random. For all $(i, j, k) \in \Omega$, let $[\mathcal{E}_0]_{ijk} = \mathcal{M}_{ijk}$, where \mathcal{M} is a tensor with independent Bernoulli ± 1 entries.

Table 2 reports the recovery results based on varying choices of the tubal rank r of \mathcal{L}_0 and the sparsity m of \mathcal{E}_0 . It can

TABLE 2: Correct recovery for random problems of varying sizes.

$r = \operatorname{rank}_t(\mathcal{L}_0) = 0.05n, m = \ \mathcal{E}_0\ _0 = 0.05n^3$						
n	r	m	$\operatorname{rank}_t(\hat{\mathcal{L}})$	$\ \hat{\mathcal{S}}\ _0$	$\frac{\ \mathcal{L} - \mathcal{L}_0\ _F}{\ \mathcal{L}_0\ _F}$	$\frac{\ \mathcal{E} - \mathcal{E}_0\ _F}{\ \mathcal{E}_0\ _F}$
100	5	5e4	5	50,029	2.6e-7	5.4e-10
200	10	4e5	10	400,234	5.9e-7	6.7e-10
$r = \operatorname{rank}_t(\mathcal{L}_0) = 0.05n, m = \ \mathcal{E}_0\ _0 = 0.1n^3$						
n	r	m	$\operatorname{rank}_t(\hat{\mathcal{L}})$	$\ \hat{\mathcal{S}}\ _0$	$\frac{\ \mathcal{L} - \mathcal{L}_0\ _F}{\ \mathcal{L}_0\ _F}$	$\frac{\ \mathcal{E} - \mathcal{E}_0\ _F}{\ \mathcal{E}_0\ _F}$
100	5	1e5	5	100,117	4.1e-7	8.2e-10
200	10	8e5	10	800,901	4.4e-7	4.5e-10
$r = \operatorname{rank}_t(\mathcal{L}_0) = 0.1n, m = \ \mathcal{E}_0\ _0 = 0.1n^3$						
n	r	m	$\operatorname{rank}_t(\hat{\mathcal{L}})$	$\ \hat{\mathcal{S}}\ _0$	$\frac{\ \mathcal{L} - \mathcal{L}_0\ _F}{\ \mathcal{L}_0\ _F}$	$\frac{\ \mathcal{E} - \mathcal{E}_0\ _F}{\ \mathcal{E}_0\ _F}$
100	10	1e5	10	101,952	4.8e-7	1.8e-9
200	20	8e5	20	815,804	4.9e-7	9.3e-10
$r = \operatorname{rank}_t(\mathcal{L}_0) = 0.1n, m = \ \mathcal{E}_0\ _0 = 0.2n^3$						
n	r	m	$\operatorname{rank}_t(\hat{\mathcal{L}})$	$\ \hat{\mathcal{S}}\ _0$	$\frac{\ \mathcal{L} - \mathcal{L}_0\ _F}{\ \mathcal{L}_0\ _F}$	$\frac{\ \mathcal{E} - \mathcal{E}_0\ _F}{\ \mathcal{E}_0\ _F}$
100	10	2e5	10	200,056	7.7e-7	4.1e-9
200	20	16e5	20	1,601,008	1.2e-6	3.1e-9



(a) TRPCA, Random Signs (b) TRPCA, Coherent Signs

Fig. 5: Correct recovery for varying tubal ranks of \mathcal{L}_0 and sparsities of \mathcal{E}_0 . Fraction of correct recoveries across 10 trials, as a function of $\operatorname{rank}_t(\mathcal{L}_0)$ (x-axis) and sparsity of \mathcal{E}_0 (y-axis). Left: $\operatorname{sgn}(\mathcal{E}_0)$ random. Right: $\mathcal{E}_0 = \mathcal{P}_\Omega \operatorname{sgn}(\mathcal{L}_0)$.

be seen that our convex program (5) gives the correct tubal rank estimation of \mathcal{L}_0 in all cases and also the relative errors $\|\hat{\mathcal{L}} - \mathcal{L}_0\|_F / \|\mathcal{L}_0\|_F$ are very small, less than 10^{-5} . The sparsity estimation of \mathcal{E}_0 is not as exact as the rank estimation, but note that the relative errors $\|\hat{\mathcal{E}} - \mathcal{E}_0\|_F / \|\mathcal{E}_0\|_F$ are all very small, less than 10^{-8} (actually much smaller than the relative errors of the recovered low-rank component). These results well verify the correct recovery phenomenon as claimed in Theorem 4.1.

5.2 Phase Transition in Tubal Rank and Sparsity

The results in Theorem 4.1 show the perfect recovery for incoherent tensor with $\operatorname{rank}_t(\mathcal{L}_0)$ on the order of $nn_3/(\mu(\log nn_3)^2)$ and the sparsity of \mathcal{E}_0 on the order of $n^2 n_3$. Now we examine the recovery phenomenon with varying tubal rank of \mathcal{L}_0 from varying sparsity of \mathcal{E}_0 . We consider the tensor \mathcal{L}_0 of size $\mathbb{R}^{n \times n \times n_3}$, where $n = 100$ and $n_3 = 50$. We generate $\mathcal{L}_0 = \mathcal{P} * \mathcal{Q}^*$, where \mathcal{P} and \mathcal{Q} are $n \times r \times n_3$ tensors with entries independently sampled from a $\mathcal{N}(0, 1/n)$ distribution. For the sparse component \mathcal{E}_0 , we consider two cases. In the first case, we assume a Bernoulli model for the support of the sparse term \mathcal{E}_0 , with random signs: each entry of \mathcal{E}_0 takes on value 0 with probability $1 - \rho$, and values ± 1 each with probability $\rho/2$. The second case chooses the support Ω in accordance with the Bernoulli model, but this time sets $\mathcal{E}_0 = \mathcal{P}_\Omega \operatorname{sgn}(\mathcal{L}_0)$. We set $\frac{r}{n} = [0.01 : 0.01 :$

2. Codes of our method available at <https://github.com/canyilu>.

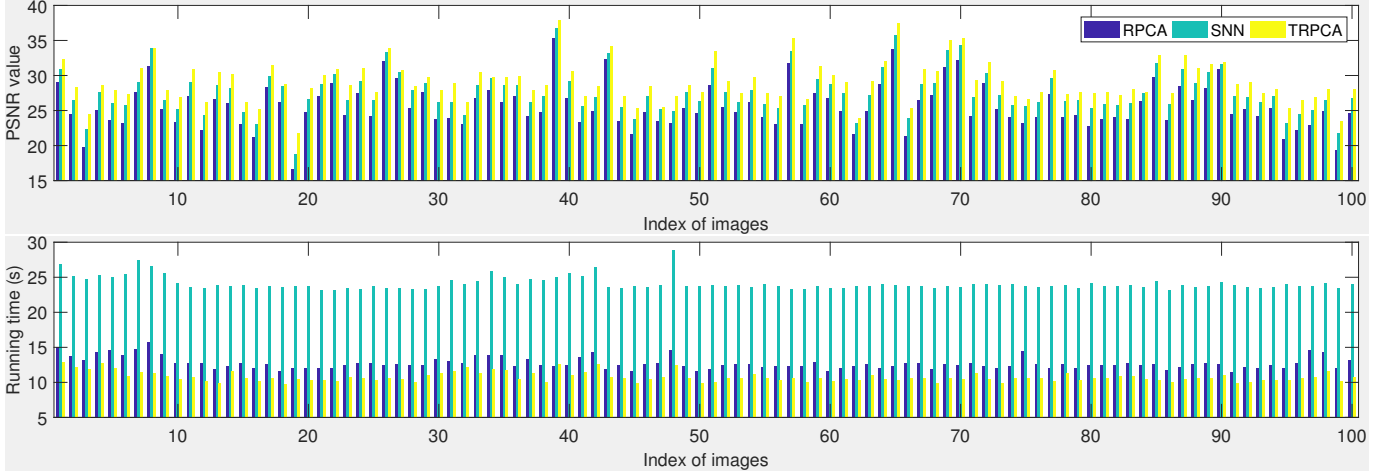


Fig. 6: Comparison of the PSNR values (top) and running time (bottom) obtained by RPCA, SNN and TRPCA on 100 images.

0.5] and $\rho_s = [0.01 : 0.01 : 0.5]$. For each $(\frac{r}{n}, \rho_s)$ -pair, we simulate 10 test instances and declare a trial to be successful if the recovered $\hat{\mathcal{L}}$ satisfies $\|\hat{\mathcal{L}} - \mathcal{L}_0\|_F / \|\mathcal{L}_0\|_F \leq 10^{-3}$. Figure 5 plots the fraction of correct recovery for each pair $(\frac{r}{n}, \rho_s)$ (black = 0% and white = 100%). It can be seen that there is a large region in which the recovery is correct in both cases. Intuitively, the experiment shows that the recovery is correct when the tubal rank of \mathcal{L}_0 is relatively low and the errors \mathcal{E}_0 is relatively sparse. Figure 5 (b) further shows that the signs of \mathcal{E}_0 are not important: recovery can be guaranteed as long as its support is chosen uniformly at random. These observations are consistent with Theorem 4.1. Similar observations can be found in the matrix RPCA case (see Figure 1 in [3]).

5.3 Application to Image Recovery

We apply TRPCA to image recovery from the corrupted images with random noises. The motivation is that the color images can be approximated by low rank matrices or tensors [18]. We will show that the recovery performance of TRPCA is still satisfactory with the suggested parameter in theory on real data.

We use 100 color images from the Berkeley Segmentation Dataset [23] for the test. The sizes of images are 321×481 or 481×321 . For each image, we randomly set 10% of pixels to random values in $[0, 255]$, and the positions of the corrupted pixels are unknown. All the 3 channels of the images are corrupted at the same positions (the corruptions are on the whole tubes). This problem is more challenging than the corruptions on 3 channels at different positions. See Figure 7 (b) for some sample images with noises. We compare our TRPCA model with RPCA [3] and SNN [13] which also own the theoretical recovery guarantee. For RPCA, we apply it on each channel separately and combine the results to obtain the recovered image. The parameter λ is set to $\lambda = 1/\sqrt{\max(n_1, n_2)}$ as suggested in theory. For SNN in (4), we find that it does not perform well when λ_i 's are set to the values suggested in theory [13]. We empirically set $\lambda = [15, 15, 1.5]$ in (4) to make SNN perform well in most cases. For our TRPCA, we format a $n_1 \times n_2$ sized image as a tensor of size $n_1 \times n_2 \times 3$. We find that such a way of tensor

construction usually performs better than some other ways. This may be due to the noises which present on the tubes. We set $\lambda = 1/\sqrt{3 \max(n_1, n_2)}$ in TRPCA. We use the Peak Signal-to-Noise Ratio (PSNR), defined as

$$\text{PSNR} = 10 \log_{10} \left(\frac{\|\mathcal{M}\|_\infty^2}{\frac{1}{n_1 n_2 n_3} \|\hat{\mathcal{X}} - \mathcal{M}\|_F^2} \right),$$

to evaluate the recovery performance.

Figure 6 gives the comparison of the PSNR values and running time on all 100 images. Some examples with the recovered images are shown in Figure 7. From these results, we have the following observations. First, both SNN and TRPCA perform much better than the matrix based RPCA. The reason is that RPCA performs on each channel independently, and thus is not able to use the information across channels. The tensor methods instead take advantage of the multi-dimensional structure of data. Second, TRPCA outperforms SNN in most cases. This not only demonstrates the superiority of our TRPCA, but also validates our recovery guarantee in Theorem 4.1 on image data. Note that SNN needs some additional effort to tune the weighted parameters λ_i 's empirically. Different from SNN which is a loose convex surrogate of the sum of Tucker rank, our TNN is a tight convex relaxation of the tensor average rank, and the recovery performance of the obtained optimal solutions has the tight recovery guarantee as RPCA. Third, we use the standard ADMM to solve RPCA, SNN and TRPCA. Figure 6 (bottom) shows that TRPCA is as efficient as RPCA, while SNN requires the highest cost in this experiment.

5.4 Application to Background Modeling

In this section, we consider the background modeling problem which aims to separate the foreground objects from the background. The frames of the background are highly correlated and thus can be modeled as a low rank tensor. The moving foreground objects occupy only a fraction of image pixels and thus can be treated as sparse errors. We solve this problem by using RPCA, SNN and TRPCA. We consider four color videos, *Hall*



(a) Original image

(b) Observed image

(c) RPCA

(d) SNN

(e) TRPCA

Index	1	2	3	4	5	6
RPCA	29.10	24.53	25.12	24.31	27.50	26.77
SNN	30.91	26.45	27.66	26.45	29.26	28.19
TRPCA	32.33	28.30	28.59	28.62	31.06	30.16

(f) Comparison of the PSNR values on the above 6 images.

Index	1	2	3	4	5	6
RPCA	14.98	13.79	14.35	12.45	12.72	15.73
SNN	26.93	25.20	25.33	23.47	23.38	28.16
TRPCA	12.96	12.24	12.76	10.70	10.64	14.31

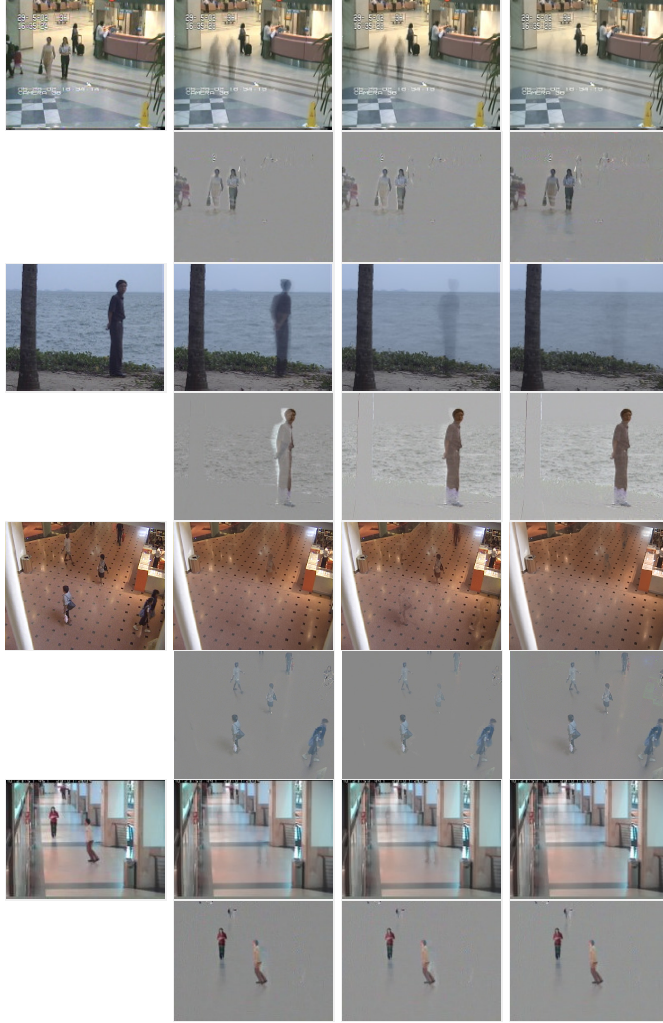
(g) Comparison of the running time (s) on the above 6 images.

Fig. 7: Recovery performance comparison on 6 example images. (a) Original image; (b) observed image; (c)-(e) recovered images by RPCA, SNN and TRPCA, respectively; (f) and (g) show the comparison of PSNR values and running time (second) on the above 6 images.

(144×176, 300), *WaterSurface* (128×160, 300), *ShoppingMall* (256×320, 100) and *ShopCorridor* (144×192, 200), where the numbers in the parentheses denote the frame size and the frame number. For each sequence with color frame size $h \times w$ and frame number k , we reshape it to a $(3hw) \times k$ matrix and use it in RPCA. To use SNN and TRPCA, we reshape the video to a $(hw) \times 3 \times k$ tensor³. The parameter of SNN in (4) is set to $\lambda = [10, 0.1, 1] \times 20$ in this experiment.

3. We observe that this way of tensor construction performs well for TRPCA, despite one has some other ways.

Figure 8 shows the performance and running time comparison of RPCA, SNN and TRPCA on the four sequences. It can be seen that the low rank components identify the main illuminations as background, while the sparse parts correspond to the motion in the scene. Generally, our TRPCA performs the best. RPCA does not perform well on the *Hall* and *WaterSurface* sequences using the default parameter. Also, TRPCA is as efficient as RPCA and SNN requires much higher computational cost. The efficiency of TRPCA is benefited from our faster way for computing tensor SVT in Algorithm 3 which is the key step for solving TRPCA.



(a) Original (b) RPCA (c) SNN (d) TRPCA

	RPCA	SNN	TRPCA
<i>Hall</i>	301.8	1553.2	323.0
<i>WaterSurface</i>	250.1	887.3	224.2
<i>ShoppingMall</i>	260.9	744.0	372.4
<i>ShopCorridor</i>	321.7	1438.6	371.3

(e) Running time (seconds) comparison

Fig. 8: Background modeling results of four surveillance video sequences. (a) Original frames; (b)-(d) low rank and sparse components obtained by RPCA, SNN and TRPCA, respectively; (e) running time comparison.

6 CONCLUSIONS AND FUTURE WORK

Based on the recently developed tensor-tensor product, which is a natural extension of the matrix-matrix product, we rigorously defined the tensor spectral norm, tensor nuclear norm and tensor average rank, such that their properties and relationships are consistent with the matrix cases. We then studied the Tensor Robust Principal Component (TRPCA) problem which aims to recover a low tubal rank tensor and a sparse tensor from their sum. We proved that under certain suitable assumptions, we can

recover both the low-rank and the sparse components exactly by simply solving a convex program whose objective is a weighted combination of the tensor nuclear norm and the ℓ_1 -norm. Benefitting from the “good” property of tensor nuclear norm, both our model and theoretical guarantee are natural extensions of RPCA. We also developed a more efficient method to compute the tensor singular value thresholding problem which is the key for solving TRPCA. Numerical experiments verify our theory and the results on images and videos demonstrate the effectiveness of our model.

There have some interesting future works. The work [7] generalizes the t-product using any invertible linear transform. With a proper choice of the invertible linear transform, it is possible to deduce a new tensor nuclear norm and solve the TRPCA problem. Beyond the convex models, the extensions to nonconvex cases are also important [21]. Finally, it is always interesting in using the developed tensor tools for real applications, e.g., image/video processing, web data analysis, and bioinformatics.

REFERENCES

- [1] K. Atkinson and W. Han. Theoretical numerical analysis: A functional analysis approach. *Texts in Applied Mathematics*, Springer, 2009.
- [2] J. Cai, E. Candès, and Z. Shen. A singular value thresholding algorithm for matrix completion. *SIAM J. Optimization*, 2010.
- [3] E. J. Candès, X. D. Li, Y. Ma, and J. Wright. Robust principal component analysis? *J. ACM*, 58(3), 2011.
- [4] Y. Chen. Incoherence-optimal matrix completion. *IEEE Trans. Information Theory*, 61(5):2909–2923, May 2015.
- [5] Anandkumar, Anima and Deng, Yuan and Ge, Rong and Mobahi, Hossein. Homotopy analysis for tensor PCA. In *Conf. on Learning Theory*, 2017.
- [6] Anandkumar, Anima and Jain, Prateek and Shi, Yang and Niranjan, Uma Naresh. Tensor vs. matrix methods: Robust tensor decomposition under block sparse perturbations. In *Artificial Intelligence and Statistics*, pages 268–276, 2016.
- [7] Kernfeld, Eric and Kilmer, Misha and Aeron, Shuchin. Tensor-tensor products with invertible linear transforms. *Linear Algebra and its Applications*, 485:545–570, 2015.
- [8] Zhang, Anru and Xia, Dong. Tensor SVD: Statistical and Computational Limits. *IEEE Trans. Information Theory*, 2018.
- [9] M. Fazel. *Matrix rank minimization with applications*. PhD thesis, Stanford University, 2002.
- [10] N. Hao, M. E. Kilmer, K. Braman, and R. C. Hoover. Facial recognition using tensor-tensor decompositions. *SIAM J. Imaging Sciences*, 6(1):437–463, 2013.
- [11] C. J. Hillar and L.-H. Lim. Most tensor problems are NP-hard. *J. ACM*, 60(6):45, 2013.
- [12] J. Hiriart-Urruty and C. Lemaréchal. *Convex Analysis and Minimization Algorithms II: Advanced Theory and Bundle Methods*. Springer, New York, 1993.
- [13] B. Huang, C. Mu, D. Goldfarb, and J. Wright. Provable models for robust low-rank tensor completion. *Pacific J. Optimization*, 11(2):339–364, 2015.
- [14] M. E. Kilmer, K. Braman, N. Hao, and R. C. Hoover. Third-order tensors as operators on matrices: A theoretical and computational framework with applications in imaging. *SIAM J. Matrix Analysis and Applications*, 34(1):148–172, 2013.
- [15] M. E. Kilmer and C. D. Martin. Factorization strategies for third-order tensors. *Linear Algebra and its Applications*, 435(3):641–658, 2011.
- [16] T. G. Kolda and B. W. Bader. Tensor decompositions and applications. *SIAM Rev.*, 51(3):455–500, 2009.
- [17] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, and Y. Ma. Robust recovery of subspace structures by low-rank representation. *IEEE Trans. Pattern Recognition and Machine Intelligence*, 2013.
- [18] J. Liu, P. Musialski, P. Wonka, and J. Ye. Tensor completion for estimating missing values in visual data. *IEEE Trans. Pattern Recognition and Machine Intelligence*, 35(1):208–220, 2013.

- [19] 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 *Proc. IEEE Conf. Computer Vision and Pattern Recognition*. IEEE, 2016.
- [20] C. Lu, J. Feng, S. Yan, and Z. Lin. A unified alternating direction method of multipliers by majorization minimization. *IEEE Trans. Pattern Recognition and Machine Intelligence*, 40(3):527–541, 2018.
- [21] C. Lu, C. Zhu, C. Xu, S. Yan, and Z. Lin. Generalized singular value thresholding. In *Proc. AAAI Conf. Artificial Intelligence*, 2015.
- [22] C. D. Martin, R. Shafer, and B. LaRue. An order- p tensor factorization with applications in imaging. *SIAM J. Scientific Computing*, 35(1):A474–A490, 2013.
- [23] D. Martin, C. Fowlkes, D. Tal, and J. Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proc. IEEE Int'l Conf. Computer Vision*, volume 2, pages 416–423. IEEE, 2001.
- [24] C. Mu, B. Huang, J. Wright, and D. Goldfarb. Square deal: Lower bounds and improved relaxations for tensor recovery. In *Proc. Int'l Conf. Machine Learning*, pages 73–81, 2014.
- [25] O. Rojo and H. Rojo. Some results on symmetric circulant matrices and on symmetric centrosymmetric matrices. *Linear algebra and its applications*, 392:211–233, 2004.
- [26] B. Romera-Paredes and M. Pontil. A new convex relaxation for tensor completion. In *Advances in Neural Information Processing Systems*, pages 2967–2975, 2013.
- [27] O. Semerci, N. Hao, M. E. Kilmer, and E. L. Miller. Tensor-based formulation and nuclear norm regularization for multienergy computed tomography. *IEEE Trans. Image Processing*, 23(4):1678–1693, 2014.
- [28] R. Tomioka, K. Hayashi, and H. Kashima. Estimation of low-rank tensors via convex optimization. *arXiv preprint arXiv:1010.0789*, 2010.
- [29] J. A. Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of Computational Mathematics*, 12(4):389–434, 2012.
- [30] R. Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- [31] A. E. Waters, A. C. Sankaranarayanan, and R. Baraniuk. SpaRCS: Recovering low-rank and sparse matrices from compressive measurements. In *Advances in Neural Information Processing Systems*, pages 1089–1097, 2011.
- [32] G. A. Watson. Characterization of the subdifferential of some matrix norms. *Linear Algebra and its Applications*, 170:33–45, 1992.
- [33] Z. Zhang and S. Aeron. Exact tensor completion using t-SVD. *IEEE Trans. Signal Processing*, 65(6):1511–1526, 2017.
- [34] Z. Zhang, G. Ely, S. Aeron, N. Hao, and M. Kilmer. Novel methods for multilinear data completion and de-noising based on tensor-SVD. In *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pages 3842–3849. IEEE, 2014.

Appendix

At the following, we give the detailed proofs of Theorem 3.1, Theorem 3.2, and the main result in Theorem 4.1. Section A first gives some notations and properties which will be used in the proofs. Section B gives the proofs of Theorem 3.1 and 3.2 in our paper. Section C provides a way for the construction of the solution to the TRPCA model, and Section D proves that the constructed solution is optimal to the TRPCA problem. Section E gives the proofs of some lemmas which are used in Section D.

APPENDIX A PRELIMINARIES

Beyond the notations introduced in the paper, we need some other notations used in the proofs. At the following, we define $\mathbf{e}_{ijk} = \hat{\mathbf{e}}_i * \hat{\mathbf{e}}_k * \hat{\mathbf{e}}_j^*$. Then we have $\mathcal{X}_{ijk} = \langle \mathcal{X}, \mathbf{e}_{ijk} \rangle$. We define the projection

$$\mathcal{P}_\Omega(\mathcal{Z}) = \sum_{ijk} \delta_{ijk} z_{ijk} \mathbf{e}_{ijk},$$

where $\delta_{ijk} = 1_{(i,j,k) \in \Omega}$, where $1_{(\cdot)}$ is the indicator function. Also Ω^c denotes the complement of Ω and $\mathcal{P}_{\Omega^\perp}$ is the projection onto Ω^c . Denote \mathcal{T} by the set

$$\mathcal{T} = \{\mathcal{U} * \mathcal{Y}^* + \mathcal{W} * \mathcal{V}^*, \mathcal{Y}, \mathcal{W} \in \mathbb{R}^{n \times r \times n_3}\}, \quad (45)$$

and by \mathcal{T}^\perp its orthogonal complement. Then the projections onto \mathcal{T} and \mathcal{T}^\perp are respectively

$$\mathcal{P}_\mathcal{T}(\mathcal{Z}) = \mathcal{U} * \mathcal{U}^* * \mathcal{Z} + \mathcal{Z} * \mathcal{V} * \mathcal{V}^* - \mathcal{U} * \mathcal{U}^* * \mathcal{Z} * \mathcal{V} * \mathcal{V}^*,$$

$$\begin{aligned} \mathcal{P}_{\mathcal{T}^\perp}(\mathcal{Z}) &= \mathcal{Z} - \mathcal{P}_\mathcal{T}(\mathcal{Z}) \\ &= (\mathcal{I}_{n_1} - \mathcal{U} * \mathcal{U}^*) * \mathcal{Z} * (\mathcal{I}_{n_2} - \mathcal{V} * \mathcal{V}^*), \end{aligned}$$

where \mathcal{I}_n denotes the $n \times n \times n_3$ identity tensor. Note that $\mathcal{P}_\mathcal{T}$ is self-adjoint. So we have

$$\begin{aligned} \|\mathcal{P}_\mathcal{T}(\mathbf{e}_{ijk})\|_F^2 &= \langle \mathcal{P}_\mathcal{T}(\mathbf{e}_{ijk}), \mathbf{e}_{ijk} \rangle \\ &= \langle \mathcal{U} * \mathcal{U}^* * \mathbf{e}_{ijk} + \mathbf{e}_{ijk} * \mathcal{V} * \mathcal{V}^*, \mathbf{e}_{ijk} \rangle \\ &\quad - \langle \mathcal{U} * \mathcal{U}^* * \mathbf{e}_{ijk} * \mathcal{V} * \mathcal{V}^*, \mathbf{e}_{ijk} \rangle \end{aligned}$$

Note that

$$\begin{aligned} &\langle \mathcal{U} * \mathcal{U}^* * \mathbf{e}_{ijk}, \mathbf{e}_{ijk} \rangle \\ &= \langle \mathcal{U} * \mathcal{U}^* * \hat{\mathbf{e}}_i * \hat{\mathbf{e}}_k * \hat{\mathbf{e}}_j^*, \hat{\mathbf{e}}_i * \hat{\mathbf{e}}_k * \hat{\mathbf{e}}_j^* \rangle \\ &= \langle \mathcal{U}^* * \hat{\mathbf{e}}_i, \mathcal{U}^* * \hat{\mathbf{e}}_i * (\hat{\mathbf{e}}_k * \hat{\mathbf{e}}_j^* * \hat{\mathbf{e}}_j * \hat{\mathbf{e}}_k^*) \rangle \\ &= \langle \mathcal{U}^* * \hat{\mathbf{e}}_i, \mathcal{U}^* * \hat{\mathbf{e}}_i \rangle \\ &= \|\mathcal{U}^* * \hat{\mathbf{e}}_i\|_F^2, \end{aligned}$$

where we use the fact that $\hat{\mathbf{e}}_k * \hat{\mathbf{e}}_j^* * \hat{\mathbf{e}}_j * \hat{\mathbf{e}}_k^* = \mathcal{I}_1$, which is the $1 \times 1 \times n_3$ identity tensor. Therefore, it is easy to see that

$$\begin{aligned} \|\mathcal{P}_\mathcal{T}(\mathbf{e}_{ijk})\|_F^2 &= \|\mathcal{U}^* * \hat{\mathbf{e}}_i\|_F^2 + \|\mathcal{V}^* * \hat{\mathbf{e}}_j\|_F^2 - \|\mathcal{U}^* * \hat{\mathbf{e}}_i * \hat{\mathbf{e}}_k * \hat{\mathbf{e}}_j^* * \mathcal{V}\|_F^2, \\ &\leq \|\mathcal{U}^* * \hat{\mathbf{e}}_i\|_F^2 + \|\mathcal{V}^* * \hat{\mathbf{e}}_j\|_F^2 \\ &\leq \frac{\mu r(n_1 + n_2)}{n_1 n_2 n_3} \end{aligned} \quad (46)$$

$$= \frac{2\mu r}{nn_3}, \text{ when } n_1 = n_2 = n. \quad (47)$$

where (46) uses the following tensor incoherence conditions

$$\max_{i=1, \dots, n_1} \|\mathcal{U}^* * \hat{\mathbf{e}}_i\|_F \leq \sqrt{\frac{\mu r}{n_1 n_3}}, \quad (48)$$

$$\max_{j=1, \dots, n_2} \|\mathcal{V}^* * \hat{\mathbf{e}}_j\|_F \leq \sqrt{\frac{\mu r}{n_2 n_3}}, \quad (49)$$

and

$$\|\mathcal{U} * \mathcal{V}^*\|_\infty \leq \sqrt{\frac{\mu r}{n_1 n_2 n_3^2}}, \quad (50)$$

which are assumed to be satisfied in Theorem 4.1 in our manuscript.

APPENDIX B

PROOFS OF THEOREM 3.1 AND THEOREM 3.2

B.1 Proof of Theorem 3.1

Proof. To complete the proof, we need the conjugate function concept. The conjugate ϕ^* of a function $\phi : C \rightarrow \mathbb{R}$, where $C \subset \mathbb{R}^n$, is defined as

$$\phi^*(\mathbf{y}) = \sup\{\langle \mathbf{y}, \mathbf{x} \rangle - \phi(\mathbf{x}) | \mathbf{x} \in C\}.$$

Note that the conjugate of the conjugate, ϕ^{**} , is the convex envelope of the function ϕ . See Theorem 1.3.5 in [12], [9]. The proofs has two steps which compute ϕ^* and ϕ^{**} , respectively.

Step 1. Computing ϕ^* . For any $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, the conjugate function of the tensor average rank

$$\phi(\mathcal{A}) = \text{rank}_a(\mathcal{A}) = \frac{1}{n_3} \text{rank}(\text{bcirc}(\mathcal{A})) = \frac{1}{n_3} \text{rank}(\bar{\mathcal{A}}),$$

on the set $S = \{\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3} | \|\mathcal{A}\| \leq 1\}$ is

$$\begin{aligned} \phi^*(\mathcal{B}) &= \sup_{\|\mathcal{A}\| \leq 1} (\langle \mathcal{B}, \mathcal{A} \rangle - \text{rank}_a(\mathcal{A})) \\ &= \sup_{\|\mathcal{A}\| \leq 1} \frac{1}{n_3} (\langle \bar{\mathcal{B}}, \bar{\mathcal{A}} \rangle - \text{rank}(\bar{\mathcal{A}})). \end{aligned}$$

Here $\bar{\mathcal{A}}, \bar{\mathcal{B}} \in \mathbb{C}^{n_1 n_3 \times n_2 n_3}$. Let $q = \min\{n_1 n_3, n_2 n_3\}$. By von Neumann's trace theorem,

$$\langle \bar{\mathcal{B}}, \bar{\mathcal{A}} \rangle \leq \sum_{i=1}^q \sigma_i(\bar{\mathcal{B}}) \sigma_i(\bar{\mathcal{A}}), \quad (51)$$

where $\sigma_i(\bar{\mathcal{A}})$ denotes the i -th largest singular value of $\bar{\mathcal{A}}$. Let $\bar{\mathcal{A}} = \bar{\mathcal{U}}_1 \bar{\mathcal{S}}_1 \bar{\mathcal{V}}_1^*$ and $\bar{\mathcal{B}} = \bar{\mathcal{U}}_2 \bar{\mathcal{S}}_2 \bar{\mathcal{V}}_2^*$ be the SVD of $\bar{\mathcal{A}}$ and $\bar{\mathcal{B}}$, respectively. Note that the equality (51) holds when

$$\bar{\mathcal{U}}_1 = \bar{\mathcal{U}}_2 \text{ and } \bar{\mathcal{V}}_1 = \bar{\mathcal{V}}_2. \quad (52)$$

So we can pick $\bar{\mathcal{U}}_1$ and $\bar{\mathcal{V}}_1$ such that (52) holds to maximize $\langle \bar{\mathcal{B}}, \bar{\mathcal{A}} \rangle$. Note that the corresponding \mathcal{U} and \mathcal{V} of $\bar{\mathcal{U}}_1$ and $\bar{\mathcal{V}}_1$ respectively are real tensors and so is \mathcal{A} in this case. Thus, we have

$$\phi^*(\mathcal{B}) = \sup_{\|\mathcal{A}\| \leq 1} \frac{1}{n_3} \left(\sum_{i=1}^q \sigma_i(\bar{\mathcal{B}}) \sigma_i(\bar{\mathcal{A}}) - \text{rank}(\bar{\mathcal{A}}) \right).$$

If $\mathcal{A} = 0$, then $\bar{\mathcal{A}} = 0$, and thus we have $\phi^*(\mathcal{B}) = 0$ for all \mathcal{B} . If $\text{rank}(\bar{\mathcal{A}}) = r$, $1 \leq r \leq q$, then $\phi^*(\mathcal{B}) = \frac{1}{n_3} (\sum_{i=1}^r \sigma_i(\bar{\mathcal{B}}) - r)$. Hence $\phi^*(\mathcal{B})$ can be expressed as

$$\begin{aligned} & n_3 \cdot \phi^*(\mathcal{B}) \\ &= \max \left\{ 0, \sigma_1(\bar{\mathcal{B}}) - 1, \dots, \sum_{i=1}^r \sigma_i(\bar{\mathcal{B}}) - r, \dots, \sum_{i=1}^q \sigma_i(\bar{\mathcal{B}}) - q \right\}. \end{aligned}$$

The largest term in this set is the one that sums all positive $(\sigma_i(\bar{\mathcal{B}}) - 1)$ terms. Thus, we have

$$\begin{aligned} & \phi^*(\mathcal{B}) \\ &= \begin{cases} 0, & \|\bar{\mathcal{B}}\| \leq 1, \\ \frac{1}{n_3} (\sum_{i=1}^r \sigma_i(\bar{\mathcal{B}}) - r), & \sigma_r(\bar{\mathcal{B}}) > 1 \text{ and } \sigma_{r+1}(\bar{\mathcal{B}}) \leq 1 \end{cases} \\ &= \frac{1}{n_3} \sum_{i=1}^q (\sigma_i(\bar{\mathcal{B}}) - 1)_+. \end{aligned}$$

Note that above $\|\bar{\mathcal{B}}\| \leq 1$ is equivalent to $\|\mathcal{B}\| \leq 1$.

Step 2. Computing ϕ^{} .** Now we compute the conjugate of ϕ^* , defined as

$$\begin{aligned} \phi^{**}(\mathcal{C}) &= \sup_{\mathcal{B}} (\langle \mathcal{C}, \mathcal{B} \rangle - \phi^*(\mathcal{B})) \\ &= \sup_{\mathcal{B}} \left(\frac{1}{n_3} \langle \bar{\mathcal{C}}, \bar{\mathcal{B}} \rangle - \phi^*(\mathcal{B}) \right), \end{aligned}$$

for all $\mathcal{C} \in S$. As before, we can choose \mathcal{B} such that

$$\phi^{**}(\mathcal{C}) = \sup_{\mathcal{B}} \left(\frac{1}{n_3} \sum_{i=1}^q \sigma_i(\bar{\mathcal{C}}) \sigma_i(\bar{\mathcal{B}}) - \phi^*(\mathcal{B}) \right).$$

At the following, we consider two cases, $\|\mathcal{C}\| > 1$ and $\|\mathcal{C}\| \leq 1$.

If $\|\mathcal{C}\| > 1$, then $\sigma_1(\bar{\mathcal{C}}) = \|\bar{\mathcal{C}}\| = \|\mathcal{C}\| > 1$. We can choose $\sigma_1(\bar{\mathcal{B}})$ large enough so that $\phi^{**}(\mathcal{C}) \rightarrow \infty$. To see this, note that in

$$\phi^{**}(\mathcal{C}) = \sup_{\mathcal{B}} \frac{1}{n_3} \left(\sum_{i=1}^q \sigma_i(\bar{\mathcal{C}}) \sigma_i(\bar{\mathcal{B}}) - \left(\sum_{i=1}^r \sigma_i(\bar{\mathcal{B}}) - r \right) \right),$$

the coefficient of $\sigma_1(\bar{\mathcal{B}})$ is $\frac{1}{n_3} (\sigma_1(\bar{\mathcal{C}}) - 1)$ which is positive.

If $\|\mathcal{C}\| \leq 1$, then $\sigma_1(\bar{\mathcal{C}}) = \|\bar{\mathcal{C}}\| = \|\mathcal{C}\| \leq 1$. If $\|\mathcal{B}\| = \|\bar{\mathcal{B}}\| \leq 1$, then $\phi^*(\mathcal{B}) = 0$ and the supremum is achieved for $\sigma_i(\bar{\mathcal{B}}) = 1$, $i = 1, \dots, q$, yielding

$$\phi^{**}(\mathcal{C}) = \frac{1}{n_3} \sum_{i=1}^q \sigma_i(\bar{\mathcal{C}}) = \frac{1}{n_3} \|\bar{\mathcal{C}}\|_* = \|\mathcal{C}\|_*.$$

If $\|\mathcal{C}\| > 1$, we show that the argument of sup is always smaller than $\|\mathcal{C}\|_*$. By adding and subtracting the term $\frac{1}{n_3} \sum_{i=1}^q \sigma_i(\bar{\mathcal{C}})$ and rearranging the terms, we have

$$\begin{aligned} & \frac{1}{n_3} \left(\sum_{i=1}^q \sigma_i(\bar{\mathcal{C}}) \sigma_i(\bar{\mathcal{B}}) - \sum_{i=1}^r (\sigma_i(\bar{\mathcal{B}}) - 1) \right) \\ &= \frac{1}{n_3} \left(\sum_{i=1}^q \sigma_i(\bar{\mathcal{C}}) \sigma_i(\bar{\mathcal{B}}) - \sum_{i=1}^r (\sigma_i(\bar{\mathcal{B}}) - 1) \right) \\ &\quad - \frac{1}{n_3} \sum_{i=1}^q \sigma_i(\bar{\mathcal{C}}) + \frac{1}{n_3} \sum_{i=1}^q \sigma_i(\bar{\mathcal{C}}) \\ &= \frac{1}{n_3} \sum_{i=1}^r (\sigma_i(\bar{\mathcal{B}}) - 1) (\sigma_i(\bar{\mathcal{C}}) - 1) \\ &\quad + \frac{1}{n_3} \sum_{i=r+1}^q (\sigma_i(\bar{\mathcal{B}}) - 1) \sigma_i(\bar{\mathcal{C}}) + \frac{1}{n_3} \sum_{i=1}^q \sigma_i(\bar{\mathcal{C}}) \\ &< \frac{1}{n_3} \sum_{i=1}^q \sigma_i(\bar{\mathcal{C}}) \\ &= \|\mathcal{C}\|_*. \end{aligned}$$

In a summary, we have shown that

$$\phi^{**}(\mathcal{C}) = \|\mathcal{C}\|_*,$$

over the set $S = \{\mathcal{C} | \|\mathcal{C}\| \leq 1\}$. Thus, $\|\mathcal{C}\|_*$ is the convex envelope of the tensor average rank $\text{rank}_a(\mathcal{C})$ over S . \square

B.2 Proof of Theorem 3.2

Proof. Let $\mathcal{G} \in \partial\|\mathcal{A}\|_*$. It is equivalent to the following statements [32]

$$\|\mathcal{A}\|_* = \langle \mathcal{G}, \mathcal{A} \rangle, \quad (53)$$

$$\|\mathcal{G}\| \leq 1. \quad (54)$$

So, to complete the proof, we only need to show that $\mathcal{G} = \mathcal{U} * \mathcal{V}^* + \mathcal{W}$, where $\mathcal{U}^* * \mathcal{W} = \mathbf{0}$, $\mathcal{W} * \mathcal{V} = \mathbf{0}$ and $\|\mathcal{W}\| \leq 1$, satisfies (53) and (54). First, we have

$$\begin{aligned} \langle \mathcal{G}, \mathcal{A} \rangle &= \langle \mathcal{U} * \mathcal{V}^* + \mathcal{W}, \mathcal{U} * \mathcal{S} * \mathcal{V}^* \rangle \\ &= \langle \mathcal{I}, \mathcal{S} \rangle + 0 \\ &= \|\mathcal{A}\|_*. \end{aligned}$$

Also, (54) is obvious when considering the property of \mathcal{W} . The proof is completed. \square

APPENDIX C DUAL CERTIFICATION

In this section, we first introduce conditions for $(\mathcal{L}_0, \mathcal{S}_0)$ to be the unique solution to TRPCA in subsection C.1. Then we construct a dual certificate in subsection C.2 which satisfies the conditions in subsection C.1, and thus our main result in Theorem 4.1 in our paper are proved.

C.1 Dual Certificates

Lemma C.1. Assume that $\|\mathcal{P}_\Omega \mathcal{P}_T\| \leq \frac{1}{2}$ and $\lambda < \frac{1}{\sqrt{n_3}}$. Then $(\mathcal{L}_0, \mathcal{S}_0)$ is the unique solution to the TRPCA problem if there is a pair $(\mathcal{W}, \mathcal{F})$ obeying

$$(\mathcal{U} * \mathcal{V}^* + \mathcal{W}) = \lambda(\text{sgn}(\mathcal{S}_0) + \mathcal{F} + \mathcal{P}_\Omega \mathcal{D}),$$

with $\mathcal{P}_T \mathcal{W} = \mathbf{0}$, $\|\mathcal{W}\| \leq \frac{1}{2}$, $\mathcal{P}_\Omega \mathcal{F} = \mathbf{0}$ and $\|\mathcal{F}\|_\infty \leq \frac{1}{2}$, and $\|\mathcal{P}_\Omega \mathcal{D}\|_F \leq \frac{1}{4}$.

Proof. For any $\mathcal{H} \neq \mathbf{0}$, $(\mathcal{L}_0 + \mathcal{H}, \mathcal{S}_0 - \mathcal{H})$ is also a feasible solution. We show that its objective is larger than that at $(\mathcal{L}_0, \mathcal{S}_0)$, hence proving that $(\mathcal{L}_0, \mathcal{S}_0)$ is the unique solution. To do this, let $\mathcal{U} * \mathcal{V}^* + \mathcal{W}_0$ be an arbitrary subgradient of the tensor nuclear norm at \mathcal{L}_0 , and $\text{sgn}(\mathcal{S}_0) + \mathcal{F}_0$ be an arbitrary subgradient of the ℓ_1 -norm at \mathcal{S}_0 . Then we have

$$\begin{aligned} &\|\mathcal{L}_0 + \mathcal{H}\|_* + \lambda\|\mathcal{S}_0 - \mathcal{H}\|_1 \\ &\geq \|\mathcal{L}_0\|_* + \lambda\|\mathcal{S}_0\|_1 + \langle \mathcal{U} * \mathcal{V}^* + \mathcal{W}_0, \mathcal{H} \rangle \\ &\quad - \lambda \langle \text{sgn}(\mathcal{S}_0) + \mathcal{F}_0, \mathcal{H} \rangle. \end{aligned}$$

Now pick \mathcal{W}_0 such that $\langle \mathcal{W}_0, \mathcal{H} \rangle = \|\mathcal{P}_{T^\perp} \mathcal{H}\|_*$ and $\langle \mathcal{F}_0, \mathcal{H} \rangle = -\|\mathcal{P}_{\Omega^\perp} \mathcal{H}\|$. We have

$$\begin{aligned} &\|\mathcal{L}_0 + \mathcal{H}\|_* + \lambda\|\mathcal{S}_0 - \mathcal{H}\|_1 \\ &\geq \|\mathcal{L}_0\|_* + \lambda\|\mathcal{S}_0\|_1 + \|\mathcal{P}_{T^\perp} \mathcal{H}\|_* + \lambda\|\mathcal{P}_{\Omega^\perp} \mathcal{H}\|_1 \\ &\quad + \langle \mathcal{U} * \mathcal{V}^* - \lambda \text{sgn}(\mathcal{S}_0), \mathcal{H} \rangle. \end{aligned}$$

By assumption

$$\begin{aligned} &|\langle \mathcal{U} * \mathcal{V}^* - \lambda \text{sgn}(\mathcal{S}_0), \mathcal{H} \rangle| \\ &\leq |\langle \mathcal{W}, \mathcal{H} \rangle| + \lambda |\langle \mathcal{F}, \mathcal{H} \rangle| + \lambda |\langle \mathcal{P}_\Omega \mathcal{D}, \mathcal{H} \rangle| \\ &\leq \beta (\|\mathcal{P}_{T^\perp} \mathcal{H}\|_* + \lambda \|\mathcal{P}_{\Omega^\perp} \mathcal{H}\|_1) + \frac{\lambda}{4} \|\mathcal{P}_\Omega \mathcal{H}\|_F, \end{aligned}$$

where $\beta = \max(\|\mathcal{W}\|, \|\mathcal{F}\|_\infty) < \frac{1}{2}$. Thus

$$\begin{aligned} &\|\mathcal{L}_0 + \mathcal{H}\|_* + \lambda\|\mathcal{S}_0 - \mathcal{H}\|_1 \\ &\geq \|\mathcal{L}_0\|_* + \lambda\|\mathcal{S}_0\|_1 + \frac{1}{2} (\|\mathcal{P}_{T^\perp} \mathcal{H}\|_* + \lambda\|\mathcal{P}_{\Omega^\perp} \mathcal{H}\|_1) \\ &\quad - \frac{\lambda}{4} \|\mathcal{P}_\Omega \mathcal{H}\|_F. \end{aligned}$$

On the other hand,

$$\begin{aligned} \|\mathcal{P}_\Omega \mathcal{H}\|_F &\leq \|\mathcal{P}_\Omega \mathcal{P}_T \mathcal{H}\|_F + \|\mathcal{P}_\Omega \mathcal{P}_{T^\perp} \mathcal{H}\|_F \\ &\leq \frac{1}{2} \|\mathcal{H}\|_F + \|\mathcal{P}_{T^\perp} \mathcal{H}\|_F \\ &\leq \frac{1}{2} \|\mathcal{P}_\Omega \mathcal{H}\|_F + \frac{1}{2} \|\mathcal{P}_{\Omega^\perp} \mathcal{H}\|_F + \|\mathcal{P}_{T^\perp} \mathcal{H}\|_F. \end{aligned}$$

Thus

$$\begin{aligned} \|\mathcal{P}_\Omega \mathcal{H}\|_F &\leq \|\mathcal{P}_{\Omega^\perp} \mathcal{H}\|_F + 2\|\mathcal{P}_{T^\perp} \mathcal{H}\|_F \\ &\leq \|\mathcal{P}_{\Omega^\perp} \mathcal{H}\|_1 + 2\sqrt{n_3} \|\mathcal{P}_{T^\perp} \mathcal{H}\|_*. \end{aligned}$$

In conclusion,

$$\begin{aligned} &\|\mathcal{L}_0 + \mathcal{H}\|_* + \lambda\|\mathcal{S}_0 - \mathcal{H}\|_1 \\ &\geq \|\mathcal{L}_0\|_* + \lambda\|\mathcal{S}_0\|_1 + \frac{1}{2} (1 - \lambda\sqrt{n_3}) \|\mathcal{P}_{T^\perp} \mathcal{H}\|_* \\ &\quad + \frac{\lambda}{4} \|\mathcal{P}_{\Omega^\perp} \mathcal{H}\|_1, \end{aligned}$$

where the last two terms are strictly positive when $\mathcal{H} \neq \mathbf{0}$. Thus, the proof is completed. \square

Lemma C.1 implies that it suffices to produce a dual certificate \mathcal{W} obeying

$$\begin{cases} \mathcal{W} \in T^\perp, \\ \|\mathcal{W}\| < \frac{1}{2}, \\ \|\mathcal{P}_\Omega (\mathcal{U} * \mathcal{V}^* + \mathcal{W} - \lambda \text{sgn}(\mathcal{S}_0))\|_F \leq \frac{\lambda}{4}, \\ \|\mathcal{P}_{\Omega^\perp} (\mathcal{U} * \mathcal{V}^* + \mathcal{W})\|_\infty < \frac{\lambda}{2}. \end{cases} \quad (55)$$

C.2 Dual Certification via the Golfing Scheme

In this subsection, we show how to construct a dual certificate obeying (55). Before we introduce our construction, our model assumes that $\Omega \sim \text{Ber}(\rho)$, or equivalently that $\Omega^c \sim \text{Ber}(1 - \rho)$. Now the distribution of Ω^c is the same as that of $\Omega^c = \Omega_1 \cup \Omega_2 \cup \dots \cup \Omega_{j_0}$, where each Ω_j follows the Bernoulli model with parameter q , which satisfies

$$\mathbb{P}((i, j, k) \in \Omega) = \mathbb{P}(\text{Bin}(j_0, q) = 0) = (1 - q)^{j_0},$$

so that the two models are the same if $\rho = (1 - q)^{j_0}$. Note that because of overlaps between the Ω_j 's, $q \geq (1 - \rho)/j_0$.

Now, we construct a dual certificate

$$\mathcal{W} = \mathcal{W}^{\mathcal{L}} + \mathcal{W}^{\mathcal{S}}, \quad (56)$$

where each component is as follows:

- 1) Construction of $\mathcal{W}^{\mathcal{L}}$ via the Golfing Scheme. Let $j_0 = 2 \log(nn_3)$ and $\Omega_j, j = 1, \dots, j_0$, be defined as previously described so that $\Omega^c = \cup_{1 \leq j \leq j_0} \Omega_j$. Then define

$$\mathcal{W}^{\mathcal{L}} = \mathcal{P}_{T^\perp} \mathcal{Y}_{j_0}, \quad (57)$$

where

$$\mathcal{Y}_j = \mathcal{Y}_{j-1} + q^{-1} \mathcal{P}_{\Omega_j} \mathcal{P}_T (\mathcal{U} * \mathcal{V}^* - \mathcal{Y}_{j-1}), \quad \mathcal{Y}_0 = \mathbf{0}.$$

- 2) Construction of $\mathcal{W}^{\mathcal{S}}$ via the Method of Least Squares. Assume that $\|\mathcal{P}_\Omega \mathcal{P}_T\| < 1/2$. Then, $\|\mathcal{P}_\Omega \mathcal{P}_T \mathcal{P}_\Omega\| < 1/4$, and thus, the operator $\mathcal{P}_\Omega - \mathcal{P}_\Omega \mathcal{P}_T \mathcal{P}_\Omega$ mapping Ω onto itself is invertible; we denote its inverse by $(\mathcal{P}_\Omega - \mathcal{P}_\Omega \mathcal{P}_T \mathcal{P}_\Omega)^{-1}$. We then set

$$\mathcal{W}^{\mathcal{S}} = \lambda \mathcal{P}_{T^\perp} (\mathcal{P}_\Omega - \mathcal{P}_\Omega \mathcal{P}_T \mathcal{P}_\Omega)^{-1} \text{sgn}(\mathcal{S}_0). \quad (58)$$

This is equivalent to

$$\mathcal{W}^{\mathcal{S}} = \lambda \mathcal{P}_{T^\perp} \sum_{k \geq 0} (\mathcal{P}_\Omega \mathcal{P}_T \mathcal{P}_\Omega)^k \text{sgn}(\mathcal{S}_0).$$

Since both $\mathcal{W}^{\mathcal{L}}$ and $\mathcal{W}^{\mathcal{S}}$ belong to T^\perp and $\mathcal{P}_\Omega \mathcal{W}^{\mathcal{S}} = \lambda \mathcal{P}_\Omega (\mathcal{I} - \mathcal{P}_T) (\mathcal{P}_\Omega - \mathcal{P}_\Omega \mathcal{P}_T \mathcal{P}_\Omega)^{-1} \text{sgn}(\mathcal{S}_0) = \lambda \text{sgn}(\mathcal{S}_0)$, we will establish that $\mathcal{W}^{\mathcal{L}} + \mathcal{W}^{\mathcal{S}}$ is a valid dual certificate if it obeys

$$\begin{cases} \|\mathcal{W}^{\mathcal{L}} + \mathcal{W}^{\mathcal{S}}\| < \frac{1}{2}, \\ \|\mathcal{P}_\Omega (\mathcal{U} * \mathcal{V}^* + \mathcal{W}^{\mathcal{L}})\|_F \leq \frac{\lambda}{4}, \\ \|\mathcal{P}_{\Omega^\perp} (\mathcal{U} * \mathcal{V}^* + \mathcal{W}^{\mathcal{L}} + \mathcal{W}^{\mathcal{S}})\|_\infty < \frac{\lambda}{2}. \end{cases} \quad (59)$$

This can be achieved by using the following two key lemmas:

Lemma C.2. Assume that $\Omega \sim \text{Ber}(\rho)$ with parameter $\rho \leq \rho_s$ for some $\rho_s > 0$. Set $j_0 = 2 \lceil \log(nn_3) \rceil$ (use $\log(n_{(1)} n_3)$ for the tensors of rectangular frontal slice). Then, the tensor $\mathcal{W}^{\mathcal{L}}$ obeys

- (a) $\|\mathcal{W}^{\mathcal{L}}\| < \frac{1}{4}$,
- (b) $\|\mathcal{P}_\Omega (\mathcal{U} * \mathcal{V}^* + \mathcal{W}^{\mathcal{L}})\|_F < \frac{\lambda}{4}$,
- (c) $\|\mathcal{P}_{\Omega^\perp} (\mathcal{U} * \mathcal{V}^* + \mathcal{W}^{\mathcal{L}})\|_\infty < \frac{\lambda}{4}$.

Lemma C.3. Assume that \mathcal{S}_0 is supported on a set Ω sampled as in Lemma C.2, and that the signs of \mathcal{S}_0 are independent and identically distributed symmetric (and independent of Ω). Then, the tensor $\mathcal{W}^{\mathcal{S}}$ (58) obeys

- (a) $\|\mathcal{W}^{\mathcal{S}}\| < \frac{1}{4}$,
- (b) $\|\mathcal{P}_{\Omega^\perp} \mathcal{W}^{\mathcal{S}}\|_\infty < \frac{\lambda}{4}$.

So the left task is to prove Lemma C.2 and Lemma C.3, which are given in Section D.

APPENDIX D PROOFS OF DUAL CERTIFICATION

This section gives the proofs of Lemma C.2 and Lemma C.3. To do this, we first introduce some lemmas with their proofs given in Section E.

Lemma D.1. For the Bernoulli sign variable $\mathcal{M} \in \mathbb{R}^{n \times n \times n_3}$ defined as

$$\mathcal{M}_{ijk} = \begin{cases} 1, & \text{w.p. } \rho/2, \\ 0, & \text{w.p. } 1 - \rho, \\ -1, & \text{w.p. } \rho/2, \end{cases} \quad (60)$$

where $\rho > 0$, there exists a function $\varphi(\rho)$ satisfying $\lim_{\rho \rightarrow 0^+} \varphi(\rho) = 0$, such that the following statement holds with large probability,

$$\|\mathcal{M}\| \leq \varphi(\rho) \sqrt{nn_3}.$$

Lemma D.2. Suppose $\Omega \sim \text{Ber}(\rho)$. Then with high probability,

$$\|\mathcal{P}_T - \rho^{-1} \mathcal{P}_T \mathcal{P}_\Omega \mathcal{P}_T\| \leq \epsilon,$$

provided that $\rho \geq C_0 \epsilon^{-2} (\mu r \log(nn_3)) / (nn_3)$ for some numerical constant $C_0 > 0$. For the tensor of rectangular frontal slice, we need $\rho \geq C_0 \epsilon^{-2} (\mu r \log(n_{(1)} n_3)) / (n_{(2)} n_3)$.

Corollary D.3. Assume that $\Omega \sim \text{Ber}(\rho)$, then $\|\mathcal{P}_\Omega \mathcal{P}_T\|^2 \leq \rho + \epsilon$, provided that $1 - \rho \geq C \epsilon^{-2} (\mu r \log(nn_3)) / (nn_3)$, where C is as in Lemma D.2. For the tensor with frontal slice, the modification is as in Lemma D.2.

Note that this corollary shows that $\|\mathcal{P}_\Omega \mathcal{P}_T\| \leq 1/2$, provided $|\Omega|$ is not too large.

Lemma D.4. Suppose that $\mathcal{Z} \in T$ is a fixed tensor, and $\Omega \sim \text{Ber}(\rho)$. Then, with high probability,

$$\|\mathcal{Z} - \rho^{-1} \mathcal{P}_T \mathcal{P}_\Omega \mathcal{Z}\|_\infty \leq \epsilon \|\mathcal{Z}\|_\infty, \quad (61)$$

provided that $\rho \geq C_0 \epsilon^{-2} (\mu r \log(nn_3)) / (nn_3)$ (for the tensor of rectangular frontal slice, $\rho \geq C_0 \epsilon^{-2} (\mu r \log(n_{(1)} n_3)) / (n_{(2)} n_3)$) for some numerical constant $C_0 > 0$.

Lemma D.5. Suppose \mathcal{Z} is fixed, and $\Omega \sim \text{Ber}(\rho)$. Then, with high probability,

$$\|(\mathcal{I} - \rho^{-1} \mathcal{P}_\Omega) \mathcal{Z}\| \leq \sqrt{\frac{C_0 nn_3 \log(nn_3)}{\rho}} \|\mathcal{Z}\|_\infty, \quad (62)$$

for some numerical constant $C_0 > 0$ provided that $\rho \geq C_0 \log(nn_3) / (nn_3)$ (or $\rho \geq C_0 \log(n_{(1)} n_3) / (n_{(2)} n_3)$ for the tensors with rectangular frontal slice).

D.1 Proof of Lemma C.2

Proof. We first introduce some notations. Set $\mathcal{Z}_j = \mathcal{U} * \mathcal{V}^* - \mathcal{P}_T \mathcal{Y}_j$ obeying

$$\mathcal{Z}_j = (\mathcal{P}_T - q^{-1} \mathcal{P}_T \mathcal{P}_{\Omega_j} \mathcal{P}_T) \mathcal{Z}_{j-1}.$$

So $\mathcal{Z}_j \in T$ for all $j \geq 0$. Also, note that when

$$q \geq C_0 \epsilon^{-2} \frac{\mu r \log(nn_3)}{nn_3}, \quad (63)$$

or for the tensors with rectangular frontal slices $q \geq C_0 \epsilon^{-2} \frac{\mu r \log(n_{(1)} n_3)}{n_{(2)} n_3}$, we have

$$\|\mathcal{Z}_j\|_\infty \leq \epsilon \|\mathcal{Z}_{j-1}\|_\infty \leq \epsilon^j \|\mathcal{U} * \mathcal{V}^*\|_\infty, \quad (64)$$

by Lemma D.4 and

$$\|\mathbf{Z}_j\|_F \leq \epsilon \|\mathbf{Z}_{j-1}\|_F \leq \epsilon^j \|\mathbf{U} * \mathbf{V}^*\|_F \leq \epsilon^j \sqrt{r}. \quad (65)$$

We assume $\epsilon \leq e^{-1}$.

1. Proof of (a). Note that $\mathbf{Y}_{j_0} = \sum_j q^{-1} \mathcal{P}_{\Omega_j} \mathbf{Z}_{j-1}$. We have

$$\begin{aligned} \|\mathbf{W}^{\mathcal{L}}\| &= \|\mathcal{P}_{T^\perp} \mathbf{Y}_{j_0}\| \leq \sum_j \|q^{-1} \mathcal{P}_{T^\perp} \mathcal{P}_{\Omega_j} \mathbf{Z}_{j-1}\| \\ &\leq \sum_j \|\mathcal{P}_{T^\perp} (q^{-1} \mathcal{P}_{\Omega_j} \mathbf{Z}_{j-1} - \mathbf{Z}_{j-1})\| \\ &\leq \sum_j \|q^{-1} \mathcal{P}_{\Omega_j} \mathbf{Z}_{j-1} - \mathbf{Z}_{j-1}\| \\ &\leq C'_0 \sqrt{\frac{nn_3 \log(nn_3)}{q}} \sum_j \|\mathbf{Z}_{j-1}\|_\infty \\ &\leq C'_0 \sqrt{\frac{nn_3 \log(nn_3)}{q}} \sum_j \epsilon^{j-1} \|\mathbf{U} * \mathbf{V}^*\|_\infty \\ &\leq C'_0 (1 - \epsilon)^{-1} \sqrt{\frac{nn_3 \log(nn_3)}{q}} \|\mathbf{U} * \mathbf{V}^*\|_\infty. \end{aligned}$$

The fourth step is from Lemma D.5 and the fifth is from (64). Now by using (63) and (50), we have

$$\|\mathbf{W}^{\mathcal{L}}\| \leq C' \epsilon,$$

for some numerical constant C' .

2. Proof of (b). Since $\mathcal{P}_{\Omega} \mathbf{Y}_{j_0} = \mathbf{0}$, $\mathcal{P}_{\Omega} (\mathbf{U} * \mathbf{V}^* + \mathcal{P}_{T^\perp} \mathbf{Y}_{j_0}) = \mathcal{P}_{\Omega} (\mathbf{U} * \mathbf{V}^* - \mathcal{P}_T \mathbf{Y}_{j_0}) = \mathcal{P}_{\Omega} (\mathbf{Z}_{j_0})$, and it follows from (65) that

$$\|\mathbf{Z}_{j_0}\|_F \leq \epsilon^{j_0} \|\mathbf{U} * \mathbf{V}^*\|_F \leq \epsilon^{j_0} \sqrt{r}.$$

Since $\epsilon \leq e^{-1}$ and $j_0 \geq 2 \log(nn_3)$, $\epsilon^{j_0} \leq (nn_3)^{-2}$ and this proves the claim.

3. Proof of (c). We have $\mathbf{U} * \mathbf{V}^* + \mathbf{W}^{\mathcal{L}} = \mathbf{Z}_{j_0} + \mathbf{Y}_{j_0}$ and know that \mathbf{Y}_{j_0} is supported on Ω^c . Therefore, since $\|\mathbf{Z}_{j_0}\|_F \leq \lambda/8$. We only need to show that $\|\mathbf{Y}_{j_0}\|_\infty \leq \lambda/8$. Indeed,

$$\begin{aligned} \|\mathbf{Y}_{j_0}\|_\infty &\leq q^{-1} \sum_j \|\mathcal{P}_{\Omega_j} \mathbf{Z}_{j-1}\|_\infty \\ &\leq q^{-1} \sum_j \|\mathbf{Z}_{j-1}\|_\infty \\ &\leq q^{-1} \sum_j \epsilon^{j-1} \|\mathbf{U} * \mathbf{V}^*\|_\infty. \end{aligned}$$

Since $\|\mathbf{U} * \mathbf{V}^*\|_\infty \leq \sqrt{\frac{\mu r}{n^2 n_3^2}}$, this gives

$$\|\mathbf{Y}_{j_0}\|_\infty \leq C' \frac{\epsilon^2}{\sqrt{\mu r (\log(nn_3))^2}},$$

for some numerical constant C' whenever q obeys (63). Since $\lambda = 1/\sqrt{nn_3}$, $\|\mathbf{Y}_{j_0}\|_\infty \leq \lambda/8$ if

$$\epsilon \leq C \left(\frac{\mu r (\log(nn_3))^2}{nn_3} \right)^{1/4}.$$

The claim is proved by using (63), (50) and sufficiently small ϵ (provided that ρ_r is sufficiently small. Note that everything is consistent since $C_0 \epsilon^{-2} \frac{\mu r \log(nn_3)}{nn_3} < 1$. \square

D.2 Proof of Lemma C.3

Proof. We denote $\mathcal{M} = \text{sgn}(\mathcal{S}_0)$ distributed as

$$\mathcal{M}_{ijk} = \begin{cases} 1, & \text{w.p. } \rho/2, \\ 0, & \text{w.p. } 1 - \rho, \\ -1, & \text{w.p. } \rho/2. \end{cases}$$

Note that for any $\sigma > 0$, $\{\|\mathcal{P}_{\Omega} \mathcal{P}_T\| \leq \sigma\}$ holds with high probability provided that ρ is sufficiently small, see Corollary D.3.

1. Proof of (a). By construction,

$$\begin{aligned} \mathbf{W}^{\mathcal{S}} &= \lambda \mathcal{P}_{T^\perp} \mathcal{M} + \lambda \mathcal{P}_{T^\perp} \sum_{k \geq 1} (\mathcal{P}_{\Omega} \mathcal{P}_T \mathcal{P}_{\Omega})^k \mathcal{M} \\ &:= \mathcal{P}_{T^\perp} \mathbf{W}_0^{\mathcal{S}} + \mathcal{P}_{T^\perp} \mathbf{W}_1^{\mathcal{S}}. \end{aligned}$$

Note that $\|\mathcal{P}_{T^\perp} \mathbf{W}_0^{\mathcal{S}}\| \leq \|\mathbf{W}_0^{\mathcal{S}}\| = \lambda \|\mathcal{M}\|$ and $\|\mathcal{P}_{T^\perp} \mathbf{W}_1^{\mathcal{S}}\| \leq \|\mathbf{W}_1^{\mathcal{S}}\| = \lambda \|\mathcal{R}(\mathcal{M})\|$, where $\mathcal{R} = \sum_{k \geq 1} (\mathcal{P}_{\Omega} \mathcal{P}_T \mathcal{P}_{\Omega})^k$. Now, we will respectively show that $\lambda \|\mathcal{M}\|$ and $\lambda \|\mathcal{R}(\mathcal{M})\|$ are small enough when ρ is sufficiently small for $\lambda = 1/\sqrt{nn_3}$. Therefor, $\|\mathbf{W}^{\mathcal{S}}\| \leq 1/4$.

1) Bound $\lambda \|\mathcal{M}\|$.

By using Lemma D.1 directly, we have that $\lambda \|\mathcal{M}\| \leq \varphi(\rho)$ is sufficiently small given $\lambda = 1/\sqrt{nn_3}$ and ρ is sufficiently small.

2) Bound $\lambda \|\mathcal{R}(\mathcal{M})\|$.

For simplicity, let $\mathbf{Z} = \mathcal{R}(\mathcal{M})$. We have

$$\|\mathbf{Z}\| = \|\bar{\mathbf{Z}}\| = \sup_{\mathbf{x} \in \mathbb{S}^{nn_3-1}} \|\bar{\mathbf{Z}} \mathbf{x}\|_2. \quad (66)$$

The optimal \mathbf{x} to (66) is an eigenvector of $\bar{\mathbf{Z}}^* \bar{\mathbf{Z}}$. Since $\bar{\mathbf{Z}}$ is a block diagonal matrix, the optimal \mathbf{x} has a block sparse structure, i.e., $\mathbf{x} \in B = \{\mathbf{x} \in \mathbb{R}^{nn_3} | \mathbf{x} = [\mathbf{x}_1^\top, \dots, \mathbf{x}_i^\top, \dots, \mathbf{x}_{n_3}^\top]^\top, \text{ with } \mathbf{x}_i \in \mathbb{R}^n, \text{ and there exists } j \text{ such that } \mathbf{x}_j \neq \mathbf{0} \text{ and } \mathbf{x}_i = \mathbf{0}, i \neq j\}$. Note that $\|\mathbf{x}\|_2 = \|\mathbf{x}_j\|_2 = 1$. Let N be the $1/2$ -net for \mathbb{S}^{n-1} of size at most 5^n (see Lemma 5.2 in [30]). Then the $1/2$ -net, denoted as N' , for B has the size at most $n_3 \cdot 5^n$. We have

$$\begin{aligned} \|\mathcal{R}(\mathcal{M})\| &= \|\text{bdiag}(\overline{\mathcal{R}(\mathcal{M})})\| \\ &= \sup_{\mathbf{x}, \mathbf{y} \in B} \langle \mathbf{x}, \text{bdiag}(\overline{\mathcal{R}(\mathcal{M})}) \mathbf{y} \rangle \\ &= \sup_{\mathbf{x}, \mathbf{y} \in B} \langle \mathbf{x} \mathbf{y}^*, \text{bdiag}(\overline{\mathcal{R}(\mathcal{M})}) \rangle \\ &= \sup_{\mathbf{x}, \mathbf{y} \in B} \langle \text{bdiag}^*(\mathbf{x} \mathbf{y}^*), \overline{\mathcal{R}(\mathcal{M})} \rangle, \end{aligned}$$

where bdiag^* , the joint operator of bdiag , maps the block diagonal matrix $\mathbf{x} \mathbf{y}^*$ to a tensor of size $n \times n \times n_3$. Let $\mathbf{Z}' = \text{bdiag}^*(\mathbf{x} \mathbf{y}^*)$ and $\mathbf{Z} = \text{ifft}(\mathbf{Z}', [], 3)$. We have

$$\begin{aligned} \|\mathcal{R}(\mathcal{M})\| &= \sup_{\mathbf{x}, \mathbf{y} \in B} \langle \mathbf{Z}', \overline{\mathcal{R}(\mathcal{M})} \rangle \\ &= \sup_{\mathbf{x}, \mathbf{y} \in B} n_3 \langle \mathbf{Z}, \mathcal{R}(\mathcal{M}) \rangle \\ &= \sup_{\mathbf{x}, \mathbf{y} \in B} n_3 \langle \mathcal{R}(\mathbf{Z}), \mathcal{M} \rangle \\ &\leq \sup_{\mathbf{x}, \mathbf{y} \in N'} 4n_3 \langle \mathcal{R}(\mathbf{Z}), \mathcal{M} \rangle. \end{aligned}$$

For a fixed pair (\mathbf{x}, \mathbf{y}) of unit-normed vectors, define the random variable

$$X(\mathbf{x}, \mathbf{y}) = \langle 4n_3 \mathcal{R}(\mathcal{Z}), \mathcal{M} \rangle.$$

Conditional on $\Omega = \text{supp}(\mathcal{M})$, the signs of \mathcal{M} are independent and identically distributed symmetric and Hoeffding's inequality gives

$$\mathbb{P}(|X(\mathbf{x}, \mathbf{y})| > t | \Omega) \leq 2 \exp \left(-\frac{2t^2}{\|4n_3 \mathcal{R}(\mathcal{Z})\|_F^2} \right).$$

Note that $\|4n_3 \mathcal{R}(\mathcal{Z})\|_F \leq 4n_3 \|\mathcal{R}\| \|\mathcal{Z}\|_F = 4\sqrt{n_3} \|\mathcal{R}\| \|\mathcal{Z}'\|_F = 4\sqrt{n_3} \|\mathcal{R}\|$. Therefore, we have

$$\mathbb{P} \left(\sup_{\mathbf{x}, \mathbf{y} \in N'} |X(\mathbf{x}, \mathbf{y})| > t | \Omega \right) \leq 2|N'|^2 \exp \left(-\frac{t^2}{8n_3 \|\mathcal{R}\|^2} \right).$$

Hence,

$$\mathbb{P}(\|\mathcal{R}(\mathcal{M})\| > t | \Omega) \leq 2|N'|^2 \exp \left(-\frac{t^2}{8n_3 \|\mathcal{R}\|^2} \right).$$

On the event $\{\|\mathcal{P}_\Omega \mathcal{P}_T\| \leq \sigma\}$,

$$\|\mathcal{R}\| \leq \sum_{k \geq 1} \sigma^{2k} = \frac{\sigma^2}{1 - \sigma^2},$$

and, therefore, unconditionally,

$$\begin{aligned} & \mathbb{P}(\|\mathcal{R}(\mathcal{M})\| > t) \\ & \leq 2|N'|^2 \exp \left(-\frac{\gamma^2 t^2}{8n_3} \right) + \mathbb{P}(\|\mathcal{P}_\Omega \mathcal{P}_T\| \geq \sigma), \quad \gamma = \frac{1 - \sigma^2}{2\sigma^2} \\ & = 2n_3^2 \cdot 5^{2n} \exp \left(-\frac{\gamma^2 t^2}{8n_3} \right) + \mathbb{P}(\|\mathcal{P}_\Omega \mathcal{P}_T\| \geq \sigma). \end{aligned}$$

Let $t = c\sqrt{nn_3}$, where c can be a small absolute constant. Then the above inequality implies that $\|\mathcal{R}(\mathcal{M})\| \leq t$ with high probability.

2. Proof of (b) Observe that

$$\mathcal{P}_\Omega^\perp \mathcal{W}^S = -\lambda \mathcal{P}_\Omega^\perp \mathcal{P}_T (\mathcal{P}_\Omega - \mathcal{P}_\Omega \mathcal{P}_T \mathcal{P}_\Omega)^{-1} \mathcal{M}.$$

Now for $(i, j, k) \in \Omega^c$, $\mathcal{W}_{ijk}^S = \langle \mathcal{W}^S, \mathbf{e}_{ijk} \rangle$, and we have $\mathcal{W}_{ijk}^S = \lambda \langle \mathcal{Q}(i, j, k), \mathcal{M} \rangle$, where $\mathcal{Q}(i, j, k)$ is the tensor $-(\mathcal{P}_\Omega - \mathcal{P}_\Omega \mathcal{P}_T \mathcal{P}_\Omega)^{-1} \mathcal{P}_\Omega \mathcal{P}_T (\mathbf{e}_{ijk})$. Conditional on $\Omega = \text{supp}(\mathcal{M})$, the signs of \mathcal{M} are independent and identically distributed symmetric, and the Hoeffding's inequality gives

$$\mathbb{P}(|\mathcal{W}_{ijk}^S| > t\lambda | \Omega) \leq 2 \exp \left(-\frac{2t^2}{\|\mathcal{Q}(i, j, k)\|_F^2} \right),$$

and

$$\begin{aligned} & \mathbb{P}(\sup_{i,j,k} |\mathcal{W}_{ijk}^S| > t\lambda/n_3 | \Omega) \\ & \leq 2n^2 n_3 \exp \left(-\frac{2t^2}{\sup_{i,j,k} \|\mathcal{Q}(i, j, k)\|_F^2} \right). \end{aligned}$$

By using (47), we have

$$\begin{aligned} \|\mathcal{P}_\Omega \mathcal{P}_T(\mathbf{e}_{ijk})\|_F & \leq \|\mathcal{P}_\Omega \mathcal{P}_T\| \|\mathcal{P}_T(\mathbf{e}_{ijk})\|_F \\ & \leq \sigma \sqrt{\frac{2\mu r}{nn_3}}, \end{aligned}$$

on the event $\{\|\mathcal{P}_\Omega \mathcal{P}_T\| \leq \sigma\}$. On the same event, we have $\|(\mathcal{P}_\Omega - \mathcal{P}_\Omega \mathcal{P}_T \mathcal{P}_\Omega)^{-1}\| \leq (1 - \sigma^2)^{-1}$ and thus $\|\mathcal{Q}(i, j, k)\|_F^2 \leq \frac{2\sigma^2}{(1 - \sigma^2)^2} \frac{\mu r}{nn_3}$. Then, unconditionally,

$$\begin{aligned} & \mathbb{P} \left(\sup_{i,j,k} |\mathcal{W}_{ijk}^S| > t\lambda \right) \\ & \leq 2n^2 n_3 \exp \left(-\frac{nn_3 \gamma^2 t^2}{\mu r} \right) + \mathbb{P}(\|\mathcal{P}_\Omega \mathcal{P}_T\| \geq \sigma), \end{aligned}$$

where $\gamma = \frac{(1 - \sigma^2)^2}{2\sigma^2}$. This proves the claim when $\mu r < \rho'_r nn_3 \log(nn_3)^{-1}$ and ρ'_r is sufficiently small. \square

APPENDIX E PROOFS OF SOME LEMMAS

Lemma E.1. [29] Consider a finite sequence $\{\mathbf{Z}_k\}$ of independent, random $n_1 \times n_2$ matrices that satisfy the assumption $\mathbb{E}\mathbf{Z}_k = \mathbf{0}$ and $\|\mathbf{Z}_k\| \leq R$ almost surely. Let $\sigma^2 = \max\{\|\sum_k \mathbb{E}[\mathbf{Z}_k \mathbf{Z}_k^*]\|, \max_k \|\sum_k \mathbb{E}[\mathbf{Z}_k^* \mathbf{Z}_k]\|\}$. Then, for any $t \geq 0$, we have

$$\begin{aligned} \mathbb{P} \left[\left\| \sum_k \mathbf{Z}_k \right\| \geq t \right] & \leq (n_1 + n_2) \exp \left(-\frac{t^2}{2\sigma^2 + \frac{2}{3}Rt} \right) \\ & \leq (n_1 + n_2) \exp \left(-\frac{3t^2}{8\sigma^2} \right), \text{ for } t \leq \frac{\sigma^2}{R}. \end{aligned}$$

Or, for any $c > 0$, we have

$$\left\| \sum_k \mathbf{Z}_k \right\| \geq 2\sqrt{c\sigma^2 \log(n_1 + n_2)} + cB \log(n_1 + n_2),$$

with probability at least $1 - (n_1 + n_2)^{1-c}$.

E.1 Proof of Lemma D.1

Proof. The proof has three steps.

Step 1: Approximation. We first introduce some notations. Let

$$\mathbf{f}_i^* \text{ be the } i\text{-th row of } \mathbf{F}_{n_3}, \text{ and } \mathbf{M}^H = \begin{bmatrix} \mathbf{M}_1^H \\ \mathbf{M}_2^H \\ \vdots \\ \mathbf{M}_n^H \end{bmatrix} \in \mathbb{R}^{nn_3 \times n}$$

be a matrix unfolded by \mathcal{M} , where $\mathbf{M}_i^H \in \mathbb{R}^{n_3 \times n}$ is the i -th horizontal slice of \mathcal{M} , i.e., $[\mathbf{M}_i^H]_{kj} = \mathcal{M}_{ikj}$. Consider that $\bar{\mathcal{M}} = \text{fft}(\mathcal{M}, [], 3)$, we have

$$\bar{\mathbf{M}}_i = \begin{bmatrix} \mathbf{f}_i^* \mathbf{M}_1^H \\ \mathbf{f}_i^* \mathbf{M}_2^H \\ \vdots \\ \mathbf{f}_i^* \mathbf{M}_n^H \end{bmatrix},$$

where $\bar{\mathbf{M}}_i \in \mathbb{R}^{n \times n}$ is the i -th frontal slice of $\bar{\mathcal{M}}$. Note that

$$\|\mathcal{M}\| = \|\bar{\mathcal{M}}\| = \max_{i=1, \dots, n_3} \|\bar{\mathbf{M}}_i\|. \quad (67)$$

Let N be the $1/2$ -net for \mathbb{S}^{n-1} of size at most 5^n (see Lemma 5.2 in [30]). Then Lemma 5.3 in [30] gives

$$\|\bar{\mathbf{M}}_i\| \leq 2 \max_{\mathbf{x} \in N} \|\bar{\mathbf{M}}_i \mathbf{x}\|_2. \quad (68)$$

So we consider to bound $\|\bar{\mathbf{M}}_i \mathbf{x}\|_2$.

Step 2: Concentration. We can express $\|\bar{\mathbf{M}}_i \mathbf{x}\|_2^2$ as a sum of independent random variables

$$\|\bar{\mathbf{M}}_i \mathbf{x}\|_2^2 = \sum_{j=1}^n (\mathbf{f}_i^* \mathbf{M}_j^H \mathbf{x})^2 := \sum_{j=1}^n z_j^2, \quad (69)$$

where $z_j = \langle \mathbf{M}_j^H, \mathbf{f}_i \mathbf{x}^* \rangle$, $j = 1, \dots, n$, are independent sub-gaussian random variables with $\mathbb{E} z_j^2 = \rho \|\mathbf{f}_i \mathbf{x}^*\|_F^2 = \rho n_3$. Using (60), we have

$$|[\mathbf{M}_j^H]_{kl}| = \begin{cases} 1, & \text{w.p. } \rho, \\ 0, & \text{w.p. } 1 - \rho. \end{cases}$$

Thus, the sub-gaussian norm of $[\mathbf{M}_j^H]_{kl}$, denoted as $\|\cdot\|_{\psi_2}$, is

$$\begin{aligned} \|[\mathbf{M}_j^H]_{kl}\|_{\psi_2} &= \sup_{p \geq 1} p^{-\frac{1}{2}} (\mathbb{E} |[\mathbf{M}_j^H]_{kl}|^p)^{\frac{1}{p}} \\ &= \sup_{p \geq 1} p^{-\frac{1}{2}} \rho^{\frac{1}{p}}. \end{aligned}$$

Define the function $\phi(x) = x^{-\frac{1}{2}} \rho^{\frac{1}{x}}$ on $[1, +\infty)$. The only stationary point occurs at $x^* = \log \rho^{-2}$. Thus,

$$\begin{aligned} \phi(x) &\leq \max(\phi(1), \phi(x^*)) \\ &= \max\left(\rho, (\log \rho^{-2})^{-\frac{1}{2}} \rho^{\frac{1}{\log \rho^{-2}}}\right) \\ &:= \psi(\rho). \end{aligned} \quad (70)$$

Therefore, $\|[\mathbf{M}_j^H]_{kl}\|_{\psi_2} \leq \psi(\rho)$. Consider that z_j is a sum of independent centered sub-gaussian random variables $[\mathbf{M}_j^H]_{kl}$'s, by using Lemma 5.9 in [30], we have $\|z_j\|_{\psi_2}^2 \leq c_1 (\psi(\rho))^2 n_3$, where c_1 is an absolute constant. Therefore, by Remark 5.18 and Lemma 5.14 in [30], $z_j^2 - \rho n_3$ are independent centered sub-exponential random variables with $\|z_j^2 - \rho n_3\|_{\psi_1} \leq 2\|z_j\|_{\psi_2}^2 \leq 4\|z_j\|_{\psi_2}^2 \leq 4c_1 (\psi(\rho))^2 n_3$.

Now, we use an exponential deviation inequality, Corollary 5.17 in [30], to control the sum (69). We have

$$\begin{aligned} &\mathbb{P}(|\|\bar{\mathbf{M}}_i \mathbf{x}\|_2^2 - \rho n n_3| \geq tn) \\ &= \mathbb{P}\left(\left|\sum_{j=1}^n (z_j^2 - \rho n_3)\right| \geq tn\right) \\ &\leq 2 \exp\left(-c_2 n \min\left(\left(\frac{t}{4c_1 (\psi(\rho))^2 n_3}\right)^2, \frac{t}{4c_1 (\psi(\rho))^2 n_3}\right)\right), \end{aligned}$$

where $c_2 > 0$. Let $t = c_3 (\psi(\rho))^2 n_3$ for some absolute constant c_3 , we have

$$\begin{aligned} &\mathbb{P}(|\|\bar{\mathbf{M}}_i \mathbf{x}\|_2^2 - \rho n n_3| \geq c_3 (\psi(\rho))^2 n n_3) \\ &\leq 2 \exp\left(-c_2 n \min\left(\left(\frac{c_3}{4c_1}\right)^2, \frac{c_3}{4c_1}\right)\right). \end{aligned}$$

Step 3 Union bound. Taking the union bound over all \mathbf{x} in the net N of cardinality $|N| \leq 5^n$, we obtain

$$\begin{aligned} &\mathbb{P}\left(\max_{\mathbf{x} \in N} |\|\bar{\mathbf{M}}_i \mathbf{x}\|_2^2 - \rho n n_3| \geq c_3 (\psi(\rho))^2 n n_3\right) \\ &\leq 2 \cdot 5^n \cdot \exp\left(-c_2 n \min\left(\left(\frac{c_3}{4c_1}\right)^2, \frac{c_3}{4c_1}\right)\right). \end{aligned}$$

Furthermore, taking the union bound over all $i = 1, \dots, n_3$, we have

$$\begin{aligned} &\mathbb{P}\left(\max_i \left|\max_{\mathbf{x} \in N} \|\bar{\mathbf{M}}_i \mathbf{x}\|_2^2 - \rho n n_3\right| \geq c_3 (\psi(\rho))^2 n n_3\right) \\ &\leq 2 \cdot 5^n \cdot n_3 \cdot \exp\left(-c_2 n \min\left(\left(\frac{c_3}{4c_1}\right)^2, \frac{c_3}{4c_1}\right)\right). \end{aligned}$$

This implies that, with high probability (when the constant c_3 is large enough),

$$\max_i \max_{\mathbf{x} \in N} \|\bar{\mathbf{M}}_i \mathbf{x}\|_2^2 \leq (\rho + c_3 (\psi(\rho))^2) n n_3. \quad (71)$$

Let $\varphi(\rho) = 2\sqrt{\rho + c_3 (\psi(\rho))^2}$ and it satisfies $\lim_{\rho \rightarrow 0^+} \varphi(\rho) = 0$ by using (70). The proof is completed by further combining (67), (68) and (71). \square

E.2 Proof of Lemma D.2

Proof. For any tensor \mathcal{Z} , we can write

$$\begin{aligned} &(\rho^{-1} \mathcal{P}_T \mathcal{P}_\Omega \mathcal{P}_T - \mathcal{P}_T) \mathcal{Z} \\ &= \sum_{ijk} (\rho^{-1} \delta_{ijk} - 1) \langle \mathbf{e}_{ijk}, \mathcal{P}_T \mathcal{Z} \rangle \mathcal{P}_T(\mathbf{e}_{ijk}) \\ &:= \sum_{ijk} \mathcal{H}_{ijk}(\mathcal{Z}) \end{aligned}$$

where $\mathcal{H}_{ijk} : \mathbb{R}^{n \times n \times n_3} \rightarrow \mathbb{R}^{n \times n \times n_3}$ is a self-adjoint random operator with $\mathbb{E}[\mathcal{H}_{ijk}] = \mathbf{0}$. Define the matrix operator $\bar{\mathbf{H}}_{ijk} : \mathbb{B} \rightarrow \mathbb{B}$, where $\mathbb{B} = \{\bar{\mathbf{B}} : \mathcal{B} \in \mathbb{R}^{n \times n \times n_3}\}$ denotes the set consists of block diagonal matrices with the blocks as the frontal slices of $\bar{\mathcal{B}}$, as

$$\bar{\mathbf{H}}_{ijk}(\bar{\mathcal{Z}}) = (\rho^{-1} \delta_{ijk} - 1) \langle \mathbf{e}_{ijk}, \mathcal{P}_T(\mathcal{Z}) \rangle \text{bdiag}(\overline{\mathcal{P}_T(\mathbf{e}_{ijk})}).$$

By the above definitions, we have $\|\mathcal{H}_{ijk}\| = \|\bar{\mathbf{H}}_{ijk}\|$ and $\|\sum_{ijk} \mathcal{H}_{ijk}\| = \|\sum_{ijk} \bar{\mathbf{H}}_{ijk}\|$. Also $\bar{\mathbf{H}}_{ijk}$ is self-adjoint and $\mathbb{E}[\bar{\mathbf{H}}_{ijk}] = \mathbf{0}$. To prove the result by the non-commutative Bernstein inequality, we need to bound $\|\bar{\mathbf{H}}_{ijk}\|$ and $\|\sum_{ijk} \mathbb{E}[\bar{\mathbf{H}}_{ijk}^2]\|$. First, we have

$$\begin{aligned} \|\bar{\mathbf{H}}_{ijk}\| &= \sup_{\|\bar{\mathcal{Z}}\|_F=1} \|\bar{\mathbf{H}}_{ijk}(\bar{\mathcal{Z}})\|_F \\ &\leq \sup_{\|\bar{\mathcal{Z}}\|_F=1} \rho^{-1} \|\mathcal{P}_T(\mathbf{e}_{ijk})\|_F \|\text{bdiag}(\overline{\mathcal{P}_T(\mathbf{e}_{ijk})})\|_F \|\mathcal{Z}\|_F \\ &= \sup_{\|\bar{\mathcal{Z}}\|_F=1} \rho^{-1} \|\mathcal{P}_T(\mathbf{e}_{ijk})\|_F^2 \|\bar{\mathcal{Z}}\|_F \\ &\leq \frac{2\mu r}{nn_3\rho}, \end{aligned}$$

where the last inequality uses (47). On the other hand, by direct computation, we have $\bar{\mathbf{H}}_{ijk}^2(\bar{\mathcal{Z}}) = (\rho^{-1} \delta_{ijk} -$

$1)^2 \langle \mathbf{e}_{ijk}, \mathbf{P}_T(\mathbf{Z}) \rangle \langle \mathbf{e}_{ijk}, \mathbf{P}_T(\mathbf{e}_{ijk}) \rangle \text{bdiag}(\overline{\mathbf{P}_T(\mathbf{e}_{ijk})})$. Note that $\mathbb{E}[(\rho^{-1}\delta_{ijk} - 1)^2] \leq \rho^{-1}$. We have

$$\begin{aligned}
& \left\| \sum_{ijk} \mathbb{E}[\bar{\mathbf{H}}_{ijk}^2(\bar{\mathbf{Z}})] \right\|_F \\
& \leq \rho^{-1} \left\| \sum_{ijk} \langle \mathbf{e}_{ijk}, \mathbf{P}_T(\mathbf{Z}) \rangle \langle \mathbf{e}_{ijk}, \mathbf{P}_T(\mathbf{e}_{ijk}) \rangle \text{bdiag}(\overline{\mathbf{P}_T(\mathbf{e}_{ijk})}) \right\|_F \\
& \leq \rho^{-1} \sqrt{n_3} \|\mathbf{P}_T(\mathbf{e}_{ijk})\|_F^2 \left\| \sum_{ijk} \langle \mathbf{e}_{ijk}, \mathbf{P}_T(\mathbf{Z}) \rangle \right\|_F \\
& = \rho^{-1} \sqrt{n_3} \|\mathbf{P}_T(\mathbf{e}_{ijk})\|_F^2 \|\mathbf{P}_T(\mathbf{Z})\|_F \\
& \leq \rho^{-1} \sqrt{n_3} \|\mathbf{P}_T(\mathbf{e}_{ijk})\|_F^2 \|\mathbf{Z}\|_F \\
& = \rho^{-1} \|\mathbf{P}_T(\mathbf{e}_{ijk})\|_F^2 \|\bar{\mathbf{Z}}\|_F \\
& \leq \frac{2\mu r}{nn_3\rho} \|\bar{\mathbf{Z}}\|_F.
\end{aligned}$$

This implies $\left\| \sum_{ijk} \mathbb{E}[\bar{\mathbf{H}}_{ijk}^2] \right\| \leq \frac{2\mu r}{nn_3\rho}$. Let $\epsilon \leq 1$. By Lemma E.1, we have

$$\begin{aligned}
& \mathbb{P}[\|\rho^{-1}\mathbf{P}_T\mathbf{P}_\Omega\mathbf{P}_T - \mathbf{P}_T\| > \epsilon] \\
& = \mathbb{P}\left[\left\| \sum_{ijk} \mathbf{H}_{ijk} \right\| > \epsilon\right] \\
& = \mathbb{P}\left[\left\| \sum_{ijk} \bar{\mathbf{H}}_{ijk} \right\| > \epsilon\right] \\
& \leq 2nn_3 \exp\left(-\frac{3}{8} \cdot \frac{\epsilon^2}{2\mu r/(nn_3\rho)}\right) \\
& \leq 2(nn_3)^{1-\frac{3}{16}C_0},
\end{aligned}$$

where the last inequality uses $\rho \geq C_0\epsilon^{-2}\mu r \log(nn_3)/(nn_3)$. Thus, $\|\rho^{-1}\mathbf{P}_T\mathbf{P}_\Omega\mathbf{P}_T - \mathbf{P}_T\| \leq \epsilon$ holds with high probability for some numerical constant C_0 . \square

E.3 Proof of Corollary D.3

Proof. From Lemma D.2, we have

$$\|\mathbf{P}_T - (1-\rho)^{-1}\mathbf{P}_T\mathbf{P}_\Omega^\perp\mathbf{P}_T\| \leq \epsilon,$$

provided that $1-\rho \geq C_0\epsilon^{-2}(\mu r \log(nn_3))/n$. Note that $\mathcal{I} = \mathbf{P}_\Omega + \mathbf{P}_\Omega^\perp$, we have

$$\|\mathbf{P}_T - (1-\rho)^{-1}\mathbf{P}_T\mathbf{P}_\Omega^\perp\mathbf{P}_T\| = (1-\rho)^{-1}(\mathbf{P}_T\mathbf{P}_\Omega\mathbf{P}_T - \rho\mathbf{P}_T).$$

Then, by the triangular inequality

$$\|\mathbf{P}_T\mathbf{P}_\Omega\mathbf{P}_T\| \leq \epsilon(1-\rho) + \rho\|\mathbf{P}_T\| = \rho + \epsilon(1-\rho).$$

The proof is completed by using $\|\mathbf{P}_\Omega\mathbf{P}_T\|^2 = \|\mathbf{P}_T\mathbf{P}_\Omega\mathbf{P}_T\|$. \square

E.4 Proof of Lemma D.4

Proof. For any tensor $\mathbf{Z} \in \mathcal{T}$, we write

$$\rho^{-1}\mathbf{P}_T\mathbf{P}_\Omega(\mathbf{Z}) = \sum_{ijk} \rho^{-1}\delta_{ijk}z_{ijk}\mathbf{P}_T(\mathbf{e}_{ijk}).$$

The (a, b, c) -th entry of $\rho^{-1}\mathbf{P}_T\mathbf{P}_\Omega(\mathbf{Z}) - \mathbf{Z}$ can be written as a sum of independent random variables, i.e.,

$$\begin{aligned}
& \langle \rho^{-1}\mathbf{P}_T\mathbf{P}_\Omega(\mathbf{Z}) - \mathbf{Z}, \mathbf{e}_{abc} \rangle \\
& = \sum_{ijk} (\rho^{-1}\delta_{ijk} - 1)z_{ijk} \langle \mathbf{P}_T(\mathbf{e}_{ijk}), \mathbf{e}_{abc} \rangle \\
& := \sum_{ijk} t_{ijk},
\end{aligned}$$

where t_{ijk} 's are independent and $\mathbb{E}(t_{ijk}) = 0$. Now we bound $|t_{ijk}|$ and $|\sum_{ijk} \mathbb{E}[t_{ijk}^2]|$. First

$$\begin{aligned}
& |t_{ijk}| \\
& \leq \rho^{-1} \|\mathbf{Z}\|_\infty \|\mathbf{P}_T(\mathbf{e}_{ijk})\|_F \|\mathbf{P}_T(\mathbf{e}_{abc})\|_F \\
& \leq \frac{2\mu r}{nn_3\rho} \|\mathbf{Z}\|_\infty.
\end{aligned}$$

Second, we have

$$\begin{aligned}
& \left| \sum_{ijk} \mathbb{E}[t_{ijk}^2] \right| \\
& \leq \rho^{-1} \|\mathbf{Z}\|_\infty^2 \sum_{ijk} \langle \mathbf{P}_T(\mathbf{e}_{ijk}), \mathbf{e}_{abc} \rangle^2 \\
& = \rho^{-1} \|\mathbf{Z}\|_\infty^2 \sum_{ijk} \langle \mathbf{e}_{ijk}, \mathbf{P}_T(\mathbf{e}_{abc}) \rangle^2 \\
& = \rho^{-1} \|\mathbf{Z}\|_\infty^2 \|\mathbf{P}_T(\mathbf{e}_{abc})\|_F^2 \\
& \leq \frac{2\mu r}{nn_3\rho} \|\mathbf{Z}\|_\infty^2.
\end{aligned}$$

Let $\epsilon \leq 1$. By Lemma E.1, we have

$$\begin{aligned}
& \mathbb{P}[\|\rho^{-1}\mathbf{P}_T\mathbf{P}_\Omega(\mathbf{Z}) - \mathbf{Z}\|_{abc} > \epsilon\|\mathbf{Z}\|_\infty] \\
& = \mathbb{P}\left[\left| \sum_{ijk} t_{ijk} \right| > \epsilon\|\mathbf{Z}\|_\infty\right] \\
& \leq 2 \exp\left(-\frac{3}{8} \cdot \frac{\epsilon^2\|\mathbf{Z}\|_\infty^2}{2\mu r\|\mathbf{Z}\|_\infty^2/(nn_3\rho)}\right) \\
& \leq 2(nn_3)^{-\frac{3}{16}C_0},
\end{aligned}$$

where the last inequality uses $\rho \geq C_0\epsilon^{-2}\mu r \log(nn_3)/(nn_3)$. Thus, $\|\rho^{-1}\mathbf{P}_T\mathbf{P}_\Omega(\mathbf{Z}) - \mathbf{Z}\|_\infty \leq \epsilon\|\mathbf{Z}\|_\infty$ holds with high probability for some numerical constant C_0 . \square

E.5 Proof of Lemma D.5

Proof. Denote the tensor $\mathbf{H}_{ijk} = (1 - \rho^{-1}\delta_{ijk})z_{ijk}\mathbf{e}_{ijk}$. Then we have

$$(\mathcal{I} - \rho^{-1}\mathbf{P}_\Omega)\mathbf{Z} = \sum_{ijk} \mathbf{H}_{ijk}.$$

Note that δ_{ijk} 's are independent random scalars. Thus, \mathbf{H}_{ijk} 's are independent random tensors and $\bar{\mathbf{H}}_{ijk}$'s are independent random

matrices. Observe that $\mathbb{E}[\bar{\mathbf{H}}_{ijk}] = \mathbf{0}$ and $\|\bar{\mathbf{H}}_{ijk}\| \leq \rho^{-1}\|\mathbf{Z}\|_\infty$. We have

$$\begin{aligned}
& \left\| \sum_{ijk} \mathbb{E}[\bar{\mathbf{H}}_{ijk}^* \bar{\mathbf{H}}_{ijk}] \right\| \\
&= \left\| \sum_{ijk} \mathbb{E}[\mathbf{H}_{ijk}^* * \mathbf{H}_{ijk}] \right\| \\
&= \left\| \sum_{ijk} \mathbb{E}[(1 - \rho^{-1}\delta_{ijk})^2] z_{ijk}^2 (\mathbf{e}_j * \mathbf{e}_j^*) \right\| \\
&= \left\| \frac{1-\rho}{\rho} \sum_{ijk} z_{ijk}^2 (\mathbf{e}_j * \mathbf{e}_j^*) \right\| \\
&\leq \frac{nn_3}{\rho} \|\mathbf{Z}\|_\infty^2.
\end{aligned}$$

A similar calculation yields $\left\| \sum_{ijk} \mathbb{E}[\bar{\mathbf{H}}_{ijk}^* \bar{\mathbf{H}}_{ijk}] \right\| \leq \rho^{-1}nn_3\|\mathbf{Z}\|_\infty^2$. Let $t = \sqrt{C_0nn_3 \log(nn_3)/\rho} \|\mathbf{Z}\|_\infty$. When $\rho \geq C_0 \log(nn_3)/(nn_3)$, we apply Lemma E.1 and obtain

$$\begin{aligned}
& \mathbb{P} [\|(\mathcal{I} - \rho^{-1}\mathcal{P}_\Omega)\mathbf{Z}\| > t] \\
&= \mathbb{P} \left[\left\| \sum_{ijk} \mathbf{H}_{ijk} \right\| > t \right] \\
&= \mathbb{P} \left[\left\| \sum_{ijk} \bar{\mathbf{H}}_{ijk} \right\| > t \right] \\
&\leq 2nn_3 \exp \left(-\frac{3}{8} \cdot \frac{C_0nn_3 \log(nn_3)\|\mathbf{Z}\|_\infty^2/\rho}{nn_3\|\mathbf{Z}\|_\infty^2/\rho} \right) \\
&\leq 2(nn_3)^{1-\frac{3}{8}C_0}.
\end{aligned}$$

Thus, $\|(\mathcal{I} - \rho^{-1}\mathcal{P}_\Omega)\mathbf{Z}\| > t$ holds with high probability for some numerical constant C_0 . \square