

Fast and Provable Tensor Robust Principal Component Analysis via Scaled Gradient Descent

Harry Dong*
CMU

Tian Tong*
CMU

Cong Ma[†]
UChicago

Yuejie Chi*
CMU

June 2022; Revised: February 2023

Abstract

An increasing number of data science and machine learning problems rely on computation with tensors, which better capture the multi-way relationships and interactions of data than matrices. When tapping into this critical advantage, a key challenge is to develop computationally efficient and provably correct algorithms for extracting useful information from tensor data that are simultaneously robust to corruptions and ill-conditioning. This paper tackles tensor robust principal component analysis (RPCA), which aims to recover a low-rank tensor from its observations contaminated by sparse corruptions, under the Tucker decomposition. To minimize the computation and memory footprints, we propose to directly recover the low-dimensional tensor factors—starting from a tailored spectral initialization—via scaled gradient descent (ScaledGD), coupled with an iteration-varying thresholding operation to adaptively remove the impact of corruptions. Theoretically, we establish that the proposed algorithm converges linearly to the true low-rank tensor at a constant rate that is independent with its condition number, as long as the level of corruptions is not too large. Empirically, we demonstrate that the proposed algorithm achieves better and more scalable performance than state-of-the-art tensor RPCA algorithms through synthetic experiments and real-world applications.

Keywords: low-rank tensors, Tucker decomposition, robust principal component analysis, scaled gradient descent, preconditioning.

Contents

1	Introduction	2
1.1	Our approach	3
1.2	Related works	4
1.3	Notation and tensor preliminaries	5
2	Main results	6
2.1	Problem formulation	6
2.2	Proposed algorithm	7
2.3	Performance guarantees	8
3	Outline of the analysis	9

*Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA 15213, USA; Emails: {harryd, ttong1, yuejiec}@andrew.cmu.edu. The work of T. Tong was completed while he was a graduate student at CMU.

[†]Department of Statistics, University of Chicago, Chicago, IL 60637, USA; Email: congma@uchicago.edu.

4	Numerical experiments	11
4.1	Experiments on synthetic data	11
4.2	Image denoising and outlier detection	12
4.3	Background subtraction via selective mode update	13
5	Conclusions	13
A	Proof of Lemma 1	18
A.1	Proof of (24)	19
A.2	Proof of (26)	22
B	Proof of Lemma 2	24
C	Proof of Lemma 3	27
D	Proof of Lemma 4	28
E	Technical lemmas	31
E.1	Tensor algebra	31
E.2	Perturbation bounds	33
E.3	Sparse outliers	36

1 Introduction

An increasing number of data science and machine learning problems rely on computation with tensors [KB09, PFS16], which better capture the multi-way relationships and interactions of data than matrices; examples include recommendation systems [KABO10], topic modeling [AGH⁺14], image processing [LMWY12], anomaly detection [LWQ⁺15], and so on. Oftentimes the data object of interest can be represented by a much smaller number of latent factors than what its ambient dimension suggests, which induces a low-rank structure in the underlying tensor. Unlike the matrix case, the flexibility of tensor modeling allows one to decompose a tensor under several choices of popular decompositions. The particular tensor decomposition studied in this paper is the Tucker decomposition, where a third-order tensor $\mathcal{X}_* \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ is *low-rank* if it can be decomposed as¹

$$\mathcal{X}_* = (\mathbf{U}_*^{(1)}, \mathbf{U}_*^{(2)}, \mathbf{U}_*^{(3)}) \cdot \mathcal{G}_*,$$

where $\mathbf{U}_*^{(1)} \in \mathbb{R}^{n_1 \times r_1}$, $\mathbf{U}_*^{(2)} \in \mathbb{R}^{n_2 \times r_2}$, $\mathbf{U}_*^{(3)} \in \mathbb{R}^{n_3 \times r_3}$ are the factor matrices along each mode, $\mathcal{G}_* \in \mathbb{R}^{r_1 \times r_2 \times r_3}$ is the core tensor, and $\{r_i\}_{i=1}^3$ are the rank of each mode; see Section 2.1 for the precise definition. If we flatten the tensor along each mode, then the obtained matrices are all correspondingly low-rank:

$$r_1 = \text{rank}(\mathcal{M}_1(\mathcal{X}_*)), \quad r_2 = \text{rank}(\mathcal{M}_2(\mathcal{X}_*)), \quad r_3 = \text{rank}(\mathcal{M}_3(\mathcal{X}_*)),$$

where $\mathcal{M}_k(\cdot)$ denotes the matricization of an input tensor along the k -th mode ($k = 1, 2, 3$). Intuitively, this means that the fibers along each mode lie in the same low-dimensional subspace. In other words, the tensor \mathcal{X}_* has a multi-linear rank $\mathbf{r} = (r_1, r_2, r_3)$, where typically $r_k \ll n_k$. Throughout the paper, we denote $n := \max_k n_k$ and $r := \max_k r_k$.

This paper tackles tensor robust principal component analysis (RPCA), which aims to recover a low-rank tensor \mathcal{X}_* from its observations contaminated by sparse corruptions. Mathematically, imagine we have access to a set of measurements given as

$$\mathcal{Y} = \mathcal{X}_* + \mathcal{S}_*,$$

where $\mathcal{S}_* \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ is a sparse tensor—in which the number of nonzero entries is much smaller than its ambient dimension—modeling corruptions or gross errors in the observations due to sensor failures,

¹Note that there are several other popular notation for denoting the Tucker decomposition; our choice is made to facilitate the presentation of the analysis.

anomalies, or adversarial perturbations. Our goal is to recover \mathcal{X}_* from the corrupted observation \mathcal{Y} in a computationally efficient and provably correct manner.

1.1 Our approach

In this paper, we propose a novel iterative method for tensor RPCA with provable convergence guarantees. To minimize the memory footprint, we aim to directly estimate the ground truth factors, collected in $\mathbf{F}_* = (\mathbf{U}_*^{(1)}, \mathbf{U}_*^{(2)}, \mathbf{U}_*^{(3)}, \mathcal{G}_*)$, via optimizing the following objective function:

$$\mathcal{L}(\mathbf{F}, \mathcal{S}) := \frac{1}{2} \left\| (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathcal{G} + \mathcal{S} - \mathcal{Y} \right\|_{\mathbf{F}}^2, \quad (1)$$

where $\mathbf{F} = (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}, \mathcal{G})$ and \mathcal{S} are the optimization variables for the tensor factors and the corruption tensor, respectively. Despite the nonconvexity of the objective function, a simple and intuitive approach is to update the tensor factors via gradient descent, which, unfortunately, converges slowly even when the problem instance is moderately ill-conditioned [HWZ20].

On a high level, our proposed method alternates between corruption pruning (i.e., updating \mathcal{S}) and factor refinements (i.e., updating \mathbf{F}). At the beginning of each iteration, we update the corruption tensor \mathcal{S} via thresholding the observation residuals as

$$\mathcal{S}_{t+1} = \mathcal{T}_{\zeta_{t+1}} \left(\mathcal{Y} - (\mathbf{U}_t^{(1)}, \mathbf{U}_t^{(2)}, \mathbf{U}_t^{(3)}) \cdot \mathcal{G}_t \right), \quad t = 0, 1, \dots \quad (2a)$$

where \mathcal{S}_{t+1} is the update of the corruption tensor at the t -th iteration, $\mathcal{T}_{\zeta_{t+1}}(\cdot)$ trims away the entries with magnitudes smaller than an iteration-varying threshold ζ_{t+1} that is carefully orchestrated, e.g., following a geometric decaying schedule. As the estimate of the data tensor $\mathcal{X}_t = (\mathbf{U}_t^{(1)}, \mathbf{U}_t^{(2)}, \mathbf{U}_t^{(3)}) \cdot \mathcal{G}_t$ gets more accurate, the observation residual becomes more aligned with the corruptions, therefore the thresholding operator (2a) becomes more effective in identifying and removing the impact of corruptions. Turning to the low-rank tensor factors \mathbf{F} , motivated by the recent success of scaled gradient descent (ScaledGD) [TMC21a, TMC21b, TMPB⁺22] for accelerating ill-conditioned low-rank estimation, we propose to update the tensor factors iteratively by descending along the scaled gradient directions:

$$\begin{aligned} \mathbf{U}_{t+1}^{(k)} &= \mathbf{U}_t^{(k)} - \eta \nabla_{\mathbf{U}_t^{(k)}} \mathcal{L}(\mathbf{F}_t, \mathcal{S}_{t+1}) (\check{\mathbf{U}}_t^{(k)\top} \check{\mathbf{U}}_t^{(k)})^{-1}, \quad k = 1, 2, 3, \quad \text{and} \\ \mathcal{G}_{t+1} &= \mathcal{G}_t - \eta \left((\mathbf{U}_t^{(1)\top} \mathbf{U}_t^{(1)})^{-1}, (\mathbf{U}_t^{(2)\top} \mathbf{U}_t^{(2)})^{-1}, (\mathbf{U}_t^{(3)\top} \mathbf{U}_t^{(3)})^{-1} \right) \cdot \nabla_{\mathcal{G}_t} \mathcal{L}(\mathbf{F}_t, \mathcal{S}_{t+1}). \end{aligned} \quad (2b)$$

Here, $\mathbf{F}_t = (\mathbf{U}_t^{(1)}, \mathbf{U}_t^{(2)}, \mathbf{U}_t^{(3)}, \mathcal{G}_t)$ is the estimate of the tensor factors at the t -th iteration, $\nabla_{\mathbf{U}^{(k)}} \mathcal{L}(\mathbf{F}, \mathcal{S})$ and $\nabla_{\mathcal{G}} \mathcal{L}(\mathbf{F}, \mathcal{S})$ are the partial derivatives of $\mathcal{L}(\mathbf{F}, \mathcal{S})$ with respect to the corresponding variables, $\eta > 0$ is the learning rate, and

$$\check{\mathbf{U}}_t^{(1)} = (\mathbf{U}_t^{(3)} \otimes \mathbf{U}_t^{(2)}) \mathcal{M}_1(\mathcal{G}_t)^\top, \quad \check{\mathbf{U}}_t^{(2)} = (\mathbf{U}_t^{(3)} \otimes \mathbf{U}_t^{(1)}) \mathcal{M}_2(\mathcal{G}_t)^\top, \quad \check{\mathbf{U}}_t^{(3)} = (\mathbf{U}_t^{(2)} \otimes \mathbf{U}_t^{(1)}) \mathcal{M}_3(\mathcal{G}_t)^\top$$

are used to construct the preconditioned directions of the gradients, with \otimes denoting the Kronecker product. With the preconditioners, ScaledGD balances the tensor factors to find better descent directions, the benefits of which are more accentuated in ill-conditioned tensors where the convergence rate of vanilla gradient descent degenerates significantly, while ScaledGD is capable of maintaining a linear rate of convergence regardless of the condition number.

Theoretical guarantees. Coupled with a tailored spectral initialization scheme, the proposed ScaledGD method converges linearly to the true low-rank tensor in both the Frobenius norm and the entrywise ℓ_∞ norm at a constant rate that is independent of its condition number, as long as the level of corruptions—measured in terms of the fraction of nonzero entries per fiber—does not exceed the order of $\frac{1}{\mu^2 \kappa r_1 r_2 r_3}$, where μ and κ are respectively the incoherence parameter and the condition number of the ground truth tensor \mathcal{X}_* (to be formally defined later). This not only enables fast global convergence by virtue of following the scaled gradients rather than the vanilla gradients [TMPB⁺22], but also lends additional robustness to finding the low-rank Tucker decomposition despite the presence of corruptions and gross errors. Moreover, our work

provides the first refined entrywise error analysis for tensor RPCA, suggesting the errors are distributed evenly across the entries when the ground low-rank truth tensor is incoherent. To corroborate the theoretical findings, we further demonstrate that the proposed ScaledGD algorithm achieves better and more scalable performance than state-of-the-art matrix and tensor RPCA algorithms through synthetic experiments and real-world applications.

Comparisons to prior art. While tensor RPCA has been previously investigated under various low-rank tensor decompositions, e.g., [LFC⁺16, AJSN16, DBBG19], the development of provably efficient algorithms under the Tucker decomposition remains scarce. The most closely related work is [CLX21], which proposed a Riemannian gradient descent algorithm for the same tensor RPCA model as ours. Their algorithm is proven to also achieve a constant rate of convergence—at a higher per-iteration expense—as long as the fraction of outliers per fiber does not exceed the order of $\min \left\{ \frac{1}{\mu_s^4 \kappa_s^{14} r^2 \log^2 n}, \frac{1}{\mu_s^{12} \kappa_s^{12} r^3} \right\}$ (cf. [CLX21, Theorem 5.1]), where μ_s and κ_s are the spikiness parameter and the worst-case condition number of \mathcal{X}_* , respectively. Using the relation $\mu \leq \mu_s^2 \kappa_s^2$ (cf. [CLX21, Lemma 13.5]) and $\kappa \leq \kappa_s$ (cf. (14)) to conservatively translate our bound, our algorithm succeeds as long as the corruption level is below the order of $\frac{1}{\mu_s^4 \kappa_s^5 r^3}$, which is still significantly higher than that allowed in [CLX21], when the outliers are evenly distributed across the fibers. See additional numerical comparisons in Section 4.

1.2 Related works

Broadly speaking, our work falls under the recent surge of developing both computationally efficient and provably correct algorithms for high-dimensional signal estimation via nonconvex optimization, which has been particularly fruitful for problems with inherent low-rank structures; we refer interested readers to the recent overviews [CLC19, CC18] for further pointers. In the rest of this section, we focus on works that are most closely related to our paper.

Provable algorithms for matrix RPCA. The matrix RPCA problem, which aims to decompose a low-rank matrix and a sparse matrix from their sum, has been heavily investigated since its introduction in the seminar papers [CLMW11, CSPW11]. Convex relaxation based approaches, which minimize a weighted sum of the nuclear norm of the data matrix and the ℓ_1 norm of the corruption matrix, have been demonstrated to achieve near-optimal performance guarantees [CLMW11, WGR⁺09, CSPW11, LCM10, CFMY21, CC14]. However, their computational and memory complexities are prohibitive when applied to large-scale problem instances; for example, solving the resulting semidefinite programs via accelerated proximal gradient descent [TY10] only results in a sublinear rate of convergence with a per-iteration complexity that scales cubically with the matrix dimension. To address the computational bottleneck, nonconvex methods have been developed to achieve both statistical and computational efficiencies simultaneously [NNS⁺14, CCW19, GWL16, YPCC16, TMC21a, CLY21]. Our tensor RPCA algorithm draws inspiration from [TMC21a, TMPB⁺22], which adopt a factored representation of the low-rank object and invoke scaled gradient updates to bypass the dependence of the convergence rate on the condition number. The matrix RPCA method in [CLY21] differs from [TMC21a] by using a threshold-based trimming procedure—which we also adopt—rather than a sorting-based one to identify the sparse matrix, for further computational savings.

Provable algorithms for tensor RPCA. Moving onto tensors, although one could unfold a tensor and feed the resulting matrices into a matrix RPCA algorithm [GQ14, ZWZM19], destroying the tensor structure through matricizations can result in suboptimal performance because it ignores the higher-order interactions [YZ16]. Therefore, it is desirable to directly operate in the tensor space. However, tensor algorithms encounter unique issues not present for matrices. For instance, while it appears straightforward to generalize the convex relaxation approach to tensors, it has been shown that computing the tensor nuclear norm is in fact NP-hard [FL18]; a similar drawback is applicable to the atomic norm formulation studied in [DBBG19]. Tensor RPCA has also been studied under different low-rank tensor decompositions, a small number of samples including the tubal rank [LFC⁺16, LFC⁺19] and the CP-rank [AJSN16, DBBG19]. These algorithms are not directly comparable with ours which uses the multilinear rank.

Robust low-rank tensor recovery. Broadly speaking, tensor RPCA concerns with reconstructing a high-dimensional tensor with certain low-dimensional structures from incomplete and corrupted observations. Pertaining to works that deal with the Tucker decomposition, [XY19] proposed a gradient descent based algorithm for tensor completion, [TMPB+22, TMC22] proposed scaled gradient descent algorithms for tensor regression and tensor completion (which our algorithm also adopts), [LZ21] proposed a Gauss-Newton algorithm for tensor regression that achieves quadratic convergence, [WCW21] proposed a Riemannian gradient method with entrywise convergence guarantees, and [ARB20] studied tensor regression assuming the underlying tensor is simultaneously low-rank and sparse.

1.3 Notation and tensor preliminaries

Throughout this paper, we use boldface calligraphic letters (e.g. \mathcal{X}) to denote tensors, and boldface capitalized letters (e.g. \mathbf{X}) to denote matrices. For any matrix \mathbf{X} , let $\sigma_i(\mathbf{X})$ be its i -th largest singular value, and $\sigma_{\max}(\mathbf{X})$ (resp. $\sigma_{\min}(\mathbf{X})$) to denote its largest (resp. smallest) nonzero singular value. Let $\|\mathbf{X}\|$, $\|\mathbf{X}\|_F$, $\|\mathbf{X}\|_{2,\infty}$, and $\|\mathbf{X}\|_\infty$ be the spectral norm, the Frobenius norm, the $\ell_{2,\infty}$ norm (largest ℓ_2 norm of the rows), and the entrywise ℓ_∞ norm of a matrix \mathbf{X} , respectively. The $r \times r$ identity matrix is denoted by \mathbf{I}_r . The set of invertible matrices in $\mathbb{R}^{r \times r}$ is denoted by $\text{GL}(r)$.

We now describe some preliminaries on tensor algebra that are used throughout this paper. For a tensor $\mathcal{X} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, let $[\mathcal{X}]_{i,j,k}$ be its (i, j, k) -th entry. For a tensor $\mathcal{X} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, suppose it can be represented via the multilinear multiplication

$$\mathcal{X} = (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathcal{G},$$

where $\mathbf{U}^{(k)} \in \mathbb{R}^{n_k \times r_k}$, $k = 1, 2, 3$, and $\mathcal{G} \in \mathbb{R}^{r_1 \times r_2 \times r_3}$. Equivalently, the entries of \mathcal{X} can be expressed as

$$[\mathcal{X}]_{i_1, i_2, i_3} = \sum_{j_1=1}^{r_1} \sum_{j_2=1}^{r_2} \sum_{j_3=1}^{r_3} [\mathbf{U}^{(1)}]_{i_1, j_1} [\mathbf{U}^{(2)}]_{i_2, j_2} [\mathbf{U}^{(3)}]_{i_3, j_3} [\mathcal{G}]_{j_1, j_2, j_3}.$$

The multilinear multiplication possesses several nice properties. A crucial one is that for any $\mathbf{B}^{(k)} \in \mathbb{R}^{r_k \times r_k}$, $k = 1, 2, 3$, it holds that

$$(\mathbf{U}^{(1)} \mathbf{B}^{(1)}, \mathbf{U}^{(2)} \mathbf{B}^{(2)}, \mathbf{U}^{(3)} \mathbf{B}^{(3)}) \cdot \mathcal{G} = (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot ((\mathbf{B}^{(1)}, \mathbf{B}^{(2)}, \mathbf{B}^{(3)}) \cdot \mathcal{G}). \quad (3)$$

In addition, if we flatten the tensor \mathcal{X} along different modes, the obtained matrices obey the following low-rank decompositions:

$$\mathcal{M}_1(\mathcal{X}) = \mathbf{U}^{(1)} \mathcal{M}_1(\mathcal{G}) (\mathbf{U}^{(3)} \otimes \mathbf{U}^{(2)})^\top = \mathbf{U}^{(1)} \check{\mathbf{U}}^{(1)\top}, \quad \check{\mathbf{U}}^{(1)} := (\mathbf{U}^{(3)} \otimes \mathbf{U}^{(2)}) \mathcal{M}_1(\mathcal{G})^\top, \quad (4a)$$

$$\mathcal{M}_2(\mathcal{X}) = \mathbf{U}^{(2)} \mathcal{M}_2(\mathcal{G}) (\mathbf{U}^{(3)} \otimes \mathbf{U}^{(1)})^\top = \mathbf{U}^{(2)} \check{\mathbf{U}}^{(2)\top}, \quad \check{\mathbf{U}}^{(2)} := (\mathbf{U}^{(3)} \otimes \mathbf{U}^{(1)}) \mathcal{M}_2(\mathcal{G})^\top, \quad (4b)$$

$$\mathcal{M}_3(\mathcal{X}) = \mathbf{U}^{(3)} \mathcal{M}_3(\mathcal{G}) (\mathbf{U}^{(2)} \otimes \mathbf{U}^{(1)})^\top = \mathbf{U}^{(3)} \check{\mathbf{U}}^{(3)\top}, \quad \check{\mathbf{U}}^{(3)} := (\mathbf{U}^{(2)} \otimes \mathbf{U}^{(1)}) \mathcal{M}_3(\mathcal{G})^\top. \quad (4c)$$

Given two tensors \mathcal{A} and \mathcal{B} , their inner product is defined as $\langle \mathcal{A}, \mathcal{B} \rangle = \sum_{i_1, i_2, i_3} \mathcal{A}_{i_1, i_2, i_3} \mathcal{B}_{i_1, i_2, i_3}$. The inner product satisfies the following property:

$$\langle (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathcal{G}, \mathcal{X} \rangle = \langle \mathcal{G}, (\mathbf{U}^{(1)\top}, \mathbf{U}^{(2)\top}, \mathbf{U}^{(3)\top}) \cdot \mathcal{X} \rangle. \quad (5)$$

Denote the Frobenius norm and the ℓ_∞ norm of \mathcal{X} as $\|\mathcal{X}\|_F = \sqrt{\langle \mathcal{X}, \mathcal{X} \rangle}$ and $\|\mathcal{X}\|_\infty = \max_{i_1, i_2, i_3} |\mathcal{X}_{i_1, i_2, i_3}|$, respectively. It follows that for $\mathbf{Q}_k \in \mathbb{R}^{r_k \times r_k}$, $k = 1, 2, 3$:

$$\|(\mathbf{Q}_1, \mathbf{Q}_2, \mathbf{Q}_3) \cdot \mathcal{G}\|_F \leq \|\mathbf{Q}_1\| \|\mathbf{Q}_2\| \|\mathbf{Q}_3\| \|\mathcal{G}\|_F. \quad (6)$$

Let $\mathbf{r} = (r_1, r_2, r_3)$. For a tensor $\mathcal{X} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, let its rank- \mathbf{r} higher-order singular value decomposition (HOSVD) $\mathcal{H}_{\mathbf{r}}(\mathcal{X})$ be

$$\mathcal{H}_{\mathbf{r}}(\mathcal{X}) = (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}, \mathcal{G}), \quad (7)$$

where $\mathbf{U}^{(k)}$ is the top r_k left singular vectors of $\mathcal{M}_k(\mathcal{X})$, $k = 1, 2, 3$, and $\mathcal{G} = (\mathbf{U}^{(1)\top}, \mathbf{U}^{(2)\top}, \mathbf{U}^{(3)\top}) \cdot \mathcal{X}$ is the core tensor. The HOSVD is an extension of the matrix SVD and can be seen as a special case of the Tucker decomposition; see [BL10] for an exposition. Although there are faster methods—such as [VVM12]—available, one straightforward way of computing the HOSVD is to obtain the singular vectors from performing matrix SVD on each matricization of \mathcal{X} . With these vectors, we can construct each $\mathbf{U}^{(k)}$, followed by finding $\mathcal{G} = (\mathbf{U}^{(1)\top}, \mathbf{U}^{(2)\top}, \mathbf{U}^{(3)\top}) \cdot \mathcal{X}$ to complete the process. In contrast to its matrix counterpart, the core tensor \mathcal{G} will not necessarily be diagonal.

2 Main results

2.1 Problem formulation

Suppose that the ground truth tensor $\mathcal{X}_\star \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ with multilinear rank $\mathbf{r} = (r_1, r_2, r_3)$ admits the following Tucker decomposition

$$\mathcal{X}_\star = (\mathbf{U}_\star^{(1)}, \mathbf{U}_\star^{(2)}, \mathbf{U}_\star^{(3)}) \cdot \mathcal{G}_\star, \quad (8)$$

where $\mathbf{U}_\star^{(1)} \in \mathbb{R}^{n_1 \times r_1}$, $\mathbf{U}_\star^{(2)} \in \mathbb{R}^{n_2 \times r_2}$, $\mathbf{U}_\star^{(3)} \in \mathbb{R}^{n_3 \times r_3}$ are the factor matrices along each mode, and $\mathcal{G}_\star \in \mathbb{R}^{r_1 \times r_2 \times r_3}$ is the core tensor. The Tucker decomposition is not unique since for any $\mathbf{Q}^{(k)} \in \text{GL}(r_k)$, $k = 1, 2, 3$, in view of (3), we have

$$(\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathcal{G} = (\mathbf{U}^{(1)}\mathbf{Q}^{(1)}, \mathbf{U}^{(2)}\mathbf{Q}^{(2)}, \mathbf{U}^{(3)}\mathbf{Q}^{(3)}) \cdot \mathcal{G}_Q$$

where $\mathcal{G}_Q = ((\mathbf{Q}^{(1)})^{-1}, (\mathbf{Q}^{(2)})^{-1}, (\mathbf{Q}^{(3)})^{-1}) \cdot \mathcal{G}$. Without loss of generality, to address ambiguity, we set the ground truth $\mathbf{F}_\star = (\mathbf{U}_\star^{(1)}, \mathbf{U}_\star^{(2)}, \mathbf{U}_\star^{(3)}, \mathcal{G}_\star)$ to satisfy that for each mode, $\mathbf{U}_\star^{(k)} \in \mathbb{R}^{n_k \times r_k}$ to have orthonormal columns, and

$$\mathcal{M}_k(\mathcal{G}_\star) \mathcal{M}_k(\mathcal{G}_\star)^\top = (\boldsymbol{\Sigma}_\star^{(k)})^2, \quad (9)$$

the squared singular value matrix $\boldsymbol{\Sigma}_\star^{(k)}$ of $\mathcal{M}_k(\mathcal{X}_\star)$. This can be easily met, for example, by taking the tensor factors \mathbf{F}_\star as the HOSVD of \mathcal{X}_\star .

Observation model and goal. Suppose that we collect a set of corrupted observations of \mathcal{X}_\star as

$$\mathcal{Y} = \mathcal{X}_\star + \mathcal{S}_\star, \quad (10)$$

where \mathcal{S}_\star is the corruption tensor. The problem of tensor RPCA seeks to separate \mathcal{X}_\star and \mathcal{S}_\star from their sum \mathcal{Y} as efficiently and accurately as possible.

Key quantities. Obviously, the tensor RPCA problem is ill-posed without imposing additional constraints on the low-rank tensor \mathcal{X}_\star and the corruption tensor \mathcal{S}_\star , which are crucial in determining the performance of the proposed algorithm. We first introduce the incoherence parameter of the tensor \mathcal{X}_\star .

Definition 1 (Incoherence). The incoherence parameter μ of \mathcal{X}_\star is defined as

$$\mu := \max_k \left\{ \frac{n_k}{r_k} \left\| \mathbf{U}_\star^{(k)} \right\|_{2,\infty}^2 \right\}, \quad (11)$$

where $\mathcal{X}_\star = (\mathbf{U}_\star^{(1)}, \mathbf{U}_\star^{(2)}, \mathbf{U}_\star^{(3)}) \cdot \mathcal{G}_\star$ is its Tucker decomposition.

The incoherence parameter roughly measures how spread the energy of \mathcal{X}_\star is over its entries—the energy is more spread as μ gets smaller. Moreover, we define a new notion of condition number that measures the conditioning of the ground truth tensor \mathcal{X}_\star as follows, which is weaker than previously used notions.

Definition 2 (Condition number). The condition number κ of \mathcal{X}_\star is defined as

$$\kappa := \frac{\min_k \sigma_{\max}(\mathcal{M}_k(\mathcal{X}_\star))}{\min_k \sigma_{\min}(\mathcal{M}_k(\mathcal{X}_\star))}. \quad (12)$$

With slight abuse of terminology, denote

$$\sigma_{\min}(\mathcal{X}_*) = \min_k \sigma_{\min}(\mathcal{M}_k(\mathcal{X}_*)) \quad (13)$$

as the minimum nonzero singular value of \mathcal{X}_* .

Remark 1. The above-defined condition number can be much smaller than the *worst-case* condition number κ_s used in prior analyses [TMPB+22, CLX21, HWZ20], which is defined as

$$\kappa_s := \frac{\max_k \sigma_{\max}(\mathcal{M}_k(\mathcal{X}_*))}{\min_k \sigma_{\min}(\mathcal{M}_k(\mathcal{X}_*))} \geq \frac{\min_k \sigma_{\max}(\mathcal{M}_k(\mathcal{X}_*))}{\min_k \sigma_{\min}(\mathcal{M}_k(\mathcal{X}_*))} = \kappa. \quad (14)$$

Furthermore, the condition number κ is also upper bounded by the largest condition number of the matricization along different modes, i.e., $\kappa \leq \max_k \kappa_k = \max_k \frac{\sigma_{\max}(\mathcal{M}_k(\mathcal{X}_*))}{\sigma_{\min}(\mathcal{M}_k(\mathcal{X}_*))}$.

Turning to the corruption tensor, we consider a deterministic sparsity model following the matrix case [CSPW11, NNS+14, YPCC16], where \mathcal{S}_* contains at most a small fraction of nonzero entries per fiber. This is captured in the following definition.

Definition 3 (α -fraction sparsity). The corruption tensor \mathcal{S}_* is α -fraction sparse, i.e., $\mathcal{S}_* \in \mathcal{S}_\alpha$, where

$$\mathcal{S}_\alpha := \left\{ \mathcal{S} \in \mathbb{R}^{n_1 \times n_2 \times n_3} : \|\mathcal{S}_{i_1, i_2, :}\|_0 \leq \alpha n_3, \|\mathcal{S}_{i_1, :, i_3}\|_0 \leq \alpha n_2, \|\mathcal{S}_{:, i_2, i_3}\|_0 \leq \alpha n_1, \right. \\ \left. \text{for all } 1 \leq i_k \leq n_k, \quad k = 1, 2, 3 \right\}. \quad (15)$$

With this setup in hand, we are now ready to describe the proposed algorithm.

2.2 Proposed algorithm

Our algorithm alternates between corruption removal and factor refinements. To remove the corruption, we use the following soft-shrinkage operator that trims the magnitudes of the entries by the amount of some carefully pre-set threshold.

Definition 4 (Soft-shrinkage operator). For an order-3 tensor \mathcal{X} , the soft-shrinkage operator $\mathcal{T}_\zeta(\cdot) : \mathbb{R}^{n_1 \times n_2 \times n_3} \mapsto \mathbb{R}^{n_1 \times n_2 \times n_3}$ with threshold $\zeta > 0$ is defined as

$$[\mathcal{T}_\zeta(\mathcal{X})]_{i_1, i_2, i_3} := \text{sgn}([\mathcal{X}]_{i_1, i_2, i_3}) \cdot \max(0, |[\mathcal{X}]_{i_1, i_2, i_3}| - \zeta).$$

The soft-shrinkage operator $\mathcal{T}_\zeta(\cdot)$ sets entries with magnitudes smaller than ζ to 0, while uniformly shrinking the magnitudes of the other entries by ζ . At the beginning of each iteration, the corruption tensor is updated via

$$\mathcal{S}_{t+1} = \mathcal{T}_{\zeta_t} \left(\mathcal{Y} - (\mathbf{U}_t^{(1)}, \mathbf{U}_t^{(2)}, \mathbf{U}_t^{(3)}) \cdot \mathcal{G}_t \right), \quad (16a)$$

with the schedule ζ_t to be specified shortly. With the newly updated estimate of the corruption tensor, the tensor factors are then updated by scaled gradient descent [TMPB+22], for which they are computed according to (2b) with respect to $\mathcal{L}(\mathbf{F}_t, \mathcal{S}_{t+1})$ in (1):

$$\mathbf{U}_{t+1}^{(k)} = \mathbf{U}_t^{(k)} - \eta \nabla_{\mathbf{U}_t^{(k)}} \mathcal{L}(\mathbf{F}_t, \mathcal{S}_{t+1}) (\check{\mathbf{U}}_t^{(k)\top} \check{\mathbf{U}}_t^{(k)})^{-1} \\ = (1 - \eta) \mathbf{U}_t^{(k)} - \eta (\mathcal{M}_k(\mathcal{S}_{t+1}) - \mathcal{M}_k(\mathcal{Y})) \check{\mathbf{U}}_t^{(k)} (\check{\mathbf{U}}_t^{(k)\top} \check{\mathbf{U}}_t^{(k)})^{-1} \quad (16b)$$

for $k = 1, 2, 3$ and

$$\mathcal{G}_{t+1} = \mathcal{G}_t - \eta \left((\mathbf{U}_t^{(1)\top} \mathbf{U}_t^{(1)})^{-1}, (\mathbf{U}_t^{(2)\top} \mathbf{U}_t^{(2)})^{-1}, (\mathbf{U}_t^{(3)\top} \mathbf{U}_t^{(3)})^{-1} \right) \cdot \nabla_{\mathcal{G}_t} \mathcal{L}(\mathbf{F}_t, \mathcal{S}_{t+1}) \\ = (1 - \eta) \mathcal{G}_t - \eta \left((\mathbf{U}_t^{(1)\top} \mathbf{U}_t^{(1)})^{-1} \mathbf{U}_t^{(1)\top}, (\mathbf{U}_t^{(2)\top} \mathbf{U}_t^{(2)})^{-1} \mathbf{U}_t^{(2)\top}, (\mathbf{U}_t^{(3)\top} \mathbf{U}_t^{(3)})^{-1} \mathbf{U}_t^{(3)\top} \right) \cdot (\mathcal{S}_{t+1} - \mathcal{Y}). \quad (16c)$$

Algorithm 1 ScaledGD for tensor robust principal component analysis

Input: the observed tensor \mathcal{Y} , the multilinear rank \mathbf{r} , learning rate η , and threshold schedule $\{\zeta_t\}_{t=0}^T$.
Initialization: $\mathcal{S}_0 = \mathcal{T}_{\zeta_0}(\mathcal{Y})$ and $(\mathbf{U}_0^{(1)}, \mathbf{U}_0^{(2)}, \mathbf{U}_0^{(3)}, \mathcal{G}_0) = \mathcal{H}_{\mathbf{r}}(\mathcal{Y} - \mathcal{S}_0)$.
for $t = 0, 1, \dots, T - 1$ **do**
 Update the corruption tensor \mathcal{S}_{t+1} via (16a);
 Update the tensor factors $\mathbf{F}_{t+1} = (\mathbf{U}_{t+1}^{(1)}, \mathbf{U}_{t+1}^{(2)}, \mathbf{U}_{t+1}^{(3)}, \mathcal{G}_{t+1})$ via (16b) and (16c);
end for
Output: the tensor factors $\mathbf{F}_T = (\mathbf{U}_T^{(1)}, \mathbf{U}_T^{(2)}, \mathbf{U}_T^{(3)}, \mathcal{G}_T)$.

Here, $\eta > 0$ is the learning rate, and

$$\check{\mathbf{U}}_t^{(1)} := (\mathbf{U}_t^{(3)} \otimes \mathbf{U}_t^{(2)}) \mathcal{M}_1(\mathcal{G}_t)^\top, \quad \check{\mathbf{U}}_t^{(2)} := (\mathbf{U}_t^{(3)} \otimes \mathbf{U}_t^{(1)}) \mathcal{M}_2(\mathcal{G}_t)^\top, \quad \text{and} \quad \check{\mathbf{U}}_t^{(3)} := (\mathbf{U}_t^{(2)} \otimes \mathbf{U}_t^{(1)}) \mathcal{M}_3(\mathcal{G}_t)^\top.$$

To complete the algorithm description, we still need to specify how to initialize the algorithm. We will estimate the tensor factors via the spectral method, by computing the HOSVD of the observation after applying the soft-shrinkage operator:

$$(\mathbf{U}_0^{(1)}, \mathbf{U}_0^{(2)}, \mathbf{U}_0^{(3)}, \mathcal{G}_0) = \mathcal{H}_{\mathbf{r}}(\mathcal{Y} - \mathcal{S}_0), \quad \text{where} \quad \mathcal{S}_0 = \mathcal{T}_{\zeta_0}(\mathcal{Y}).$$

Altogether, we arrive at Algorithm 1, which we still dub as ScaledGD for simplicity.

Computational benefits. It is worth highlighting that the proposed tensor RPCA algorithm possesses several computational benefits which might be of interest in applications.

- *Advantages over matrix RPCA algorithms.* While it is possible to matricize the input tensor and then apply the matrix RPCA algorithms, they can only exploit the low-rank structure along the mode that the tensor is unfolded, rather than along multiple rows simultaneously as in the tensor RPCA algorithm. In addition, the space complexity of storing and computing the factors is much higher for the matrix RPCA algorithms, where the size of the factors become multiplicative in terms of the tensor dimensions due to unfolding, rather than linear as in the tensor RPCA algorithm.
- *Generalization to N -th order tensors.* Although the description of Algorithm 1 is tailored to an order-3 tensor, our algorithm is easily generalizable to any N -th order tensor; in fact, Algorithm 1 can be applied almost verbatim by redefining

$$\check{\mathbf{U}}_t^{(k)} = (\mathbf{U}_t^{(N)} \otimes \dots \otimes \mathbf{U}_t^{(k+1)} \otimes \mathbf{U}_t^{(k-1)} \otimes \dots \otimes \mathbf{U}_t^{(1)}) \mathcal{M}_k(\mathcal{G}_t)^\top, \quad k = 1, \dots, N$$

to its natural high-order counterpart. This extension is numerically evaluated in our experiments in Section 4.

- *Parallelizability.* At each iteration of the proposed algorithm, each tensor factor is updated independently as done in (16b) and (16c), therefore we can update them in a parallel manner. This improvement becomes more apparent as the order of the tensor increases.
- *Selective modes to update:* If we know the underlying ground truth tensor is only low-rank along certain mode, we can choose to skip the iterative updates of the rest of the modes after initialization to reduce computational costs, which we demonstrate empirically in Section 4.3.

2.3 Performance guarantees

Motivated by the analysis in [TMPB⁺22], we consider the following distance metric, which not only resolves the ambiguity in the Tucker decomposition, but also takes the preconditioning factor into consideration.

Definition 5 (Distance metric). Letting $\mathbf{F} := (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}, \mathcal{G})$ and $\mathbf{F}_\star := (\mathbf{U}_\star^{(1)}, \mathbf{U}_\star^{(2)}, \mathbf{U}_\star^{(3)}, \mathcal{G}_\star)$, denote

$$\text{dist}^2(\mathbf{F}, \mathbf{F}_\star) := \inf_{\mathbf{Q}^{(k)} \in \text{GL}(r_k)} \sum_{k=1}^3 \left\| (\mathbf{U}^{(k)} \mathbf{Q}^{(k)} - \mathbf{U}_\star^{(k)}) \Sigma_\star^{(k)} \right\|_{\text{F}}^2 + \left\| ((\mathbf{Q}^{(1)})^{-1}, (\mathbf{Q}^{(2)})^{-1}, (\mathbf{Q}^{(3)})^{-1}) \cdot \mathcal{G} - \mathcal{G}_\star \right\|_{\text{F}}^2, \quad (17)$$

where we recall $\Sigma_\star^{(k)}$ is the singular value matrix of $\mathcal{M}_k(\mathcal{X}_\star)$, $k = 1, 2, 3$. Moreover, if the infimum is attained at the arguments $\{\mathbf{Q}^{(k)}\}_{k=1}^3$, they are called the optimal alignment matrices between \mathbf{F} and \mathbf{F}_\star .

Fortunately, the proposed ScaledGD algorithm (cf. Algorithm 1) provably recovers the ground truth tensor—as long as the fraction of corruptions is not too large—with proper choices of the tuning parameters, as captured in following theorem.

Theorem 1. Let $\mathcal{Y} = \mathcal{X}_\star + \mathcal{S}_\star \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, where \mathcal{X}_\star is μ -incoherent with multilinear rank $\mathbf{r} = (r_1, r_2, r_3)$, and \mathcal{S}_\star is α -sparse. Suppose that the threshold values $\{\zeta_k\}_{k=0}^\infty$ obey that $\|\mathcal{X}_\star\|_\infty \leq \zeta_0 \leq 2\|\mathcal{X}_\star\|_\infty$ and $\zeta_{t+1} = \rho\zeta_t$, $t \geq 1$, for some properly tuned $\zeta_1 := 8\sqrt{\frac{\mu^3 r_1 r_2 r_3}{n_1 n_2 n_3}} \sigma_{\min}(\mathcal{X}_\star)$ and $\frac{1}{7} \leq \eta \leq \frac{1}{4}$, where $\rho = 1 - 0.45\eta$. Then, the iterates $\mathcal{X}_t = (\mathbf{U}_t^{(1)}, \mathbf{U}_t^{(2)}, \mathbf{U}_t^{(3)}) \cdot \mathcal{G}_t$ satisfy

$$\|\mathcal{X}_t - \mathcal{X}_\star\|_{\text{F}} \leq 0.03\rho^t \sigma_{\min}(\mathcal{X}_\star), \quad (18a)$$

$$\|\mathcal{X}_t - \mathcal{X}_\star\|_\infty \leq 8\rho^t \sqrt{\frac{\mu^3 r_1 r_2 r_3}{n_1 n_2 n_3}} \sigma_{\min}(\mathcal{X}_\star), \quad (18b)$$

$$\|\mathcal{S}_t - \mathcal{S}_\star\|_\infty \leq 16\rho^{t-1} \sqrt{\frac{\mu^3 r_1 r_2 r_3}{n_1 n_2 n_3}} \sigma_{\min}(\mathcal{X}_\star) \quad (18c)$$

for all $t \geq 0$, as long as the level of corruptions obeys $\alpha \leq \frac{c_0}{\mu^2 r_1 r_2 r_3 \kappa}$ for some sufficiently small $c_0 > 0$.

The value of ρ was selected to simplify the proof and should not be taken as an optimal convergence rate. In a nutshell, Theorem 1 has the following immediate consequences:

- **Exact recovery.** Upon appropriate choices of the parameters, if the level of corruptions α is small enough, i.e. not exceeding the order of $\frac{1}{\mu^2 r_1 r_2 r_3 \kappa}$, we can ensure that the proposed Algorithm 1 exactly recovers the ground truth tensor \mathcal{X}_\star even when the gross corruptions are arbitrary and adversarial. As mentioned earlier, our result significantly enlarges the range of allowable corruption levels for exact recovery when the outliers are evenly distributed across the fibers, compared with the prior art established in [CLX21].
- **Constant linear rate of convergence.** The proposed ScaledGD algorithm (cf. Algorithm 1) finds the ground truth tensor at a *constant* linear rate, which is independent of the condition number, from a carefully designed spectral initialization. Consequently, the proposed ScaledGD algorithm inherits the computational robustness against ill-conditioning as [TMPB⁺22], even in the presence of gross outliers, as long as the thresholding operations are properly carried out.
- **Refined entrywise error guarantees.** Furthermore, when $\mu = O(1)$ and $r = O(1)$, the entrywise error bound (18b)—which is smaller than the Frobenius error (18a) by a factor of $\sqrt{\frac{1}{n_1 n_2 n_3}}$ —suggests the errors are distributed in an evenly manner across the entries for incoherent and low-rank tensors. The same applies to the entrywise error bound of the sparse tensor (18c) which exhibits similar behavior as (18b). To the best of our knowledge, this is the first time such a refined entrywise error analysis is established for tensor RPCA.

3 Outline of the analysis

In this section, we outline the proof of Theorem 1. The proof is inductive in nature, where we aim to establish the following induction hypothesis at all the iterations:

$$\text{dist}(\mathbf{F}_t, \mathbf{F}_\star) \leq \epsilon_0 \rho^t \sigma_{\min}(\mathcal{X}_\star), \quad (19a)$$

$$\max_k \left\{ \sqrt{\frac{n_k}{r_k}} \left\| (\mathbf{U}_t^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)}) \Sigma_\star^{(k)} \right\|_{2,\infty} \right\} \leq \rho^t \sqrt{\mu} \sigma_{\min}(\mathcal{X}_\star), \quad (19b)$$

where $\rho = 1 - 0.45\eta$, $\epsilon_0 < 0.01$ is some sufficiently small constant, and $\{\mathbf{Q}_t^{(k)}\}_{k=1}^3$ are the optimal alignment matrices between \mathbf{F}_t and \mathbf{F}_\star . The claims (18) in Theorem 1 follow immediately with the aid of Lemma 10 and Lemma 8 (see Appendix E). The following set of lemmas, whose proofs are deferred to Appendix A, establishes the induction hypothesis (19) for both the induction case and the base case.

Induction: local contraction. We start by outlining the local contraction of the proposed Algorithm 1, by establishing the induction hypothesis (19) continues to hold at the $(t+1)$ -th iteration, assuming it holds at the t -th iteration, as long as the corruption level is not too large.

Lemma 1 (Distance contraction). *Let $\mathbf{Y} = \mathcal{X}_\star + \mathcal{S}_\star \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, where \mathcal{X}_\star is μ -incoherent with multilinear rank $\mathbf{r} = (r_1, r_2, r_3)$, and \mathcal{S}_\star is α -sparse. Let $\mathbf{F}_t := (\mathbf{U}_t^{(1)}, \mathbf{U}_t^{(2)}, \mathbf{U}_t^{(3)}, \mathbf{G}_t)$ be the t -th iterate of Algorithm 1. Suppose that the induction hypothesis (19) holds at the t -th iteration. Under the assumption $\alpha \leq \frac{c_0 \epsilon_0}{\sqrt{\mu^3 r_1 r_2 r_3 r}}$ for some sufficiently small constant c_0 and the choice of ζ_{t+1} in Theorem 1, the $(t+1)$ -th iterate \mathbf{F}_{t+1} satisfies*

$$\text{dist}(\mathbf{F}_{t+1}, \mathbf{F}_\star) \leq \epsilon_0 \rho^{t+1} \sigma_{\min}(\mathcal{X}_\star)$$

as long as $\eta \leq 1/4$.

While Lemma 1 guarantees the contraction of the distance metric, the next Lemma 2 establishes the contraction of the incoherence metric, so that we can repeatedly apply Lemma 1 and Lemma 2 for induction.

Lemma 2 (Incoherence contraction). *Let $\mathbf{Y} = \mathcal{X}_\star + \mathcal{S}_\star \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, where \mathcal{X}_\star is μ -incoherent with multilinear rank $\mathbf{r} = (r_1, r_2, r_3)$, and \mathcal{S}_\star is α -sparse. Let $\mathbf{F}_t := (\mathbf{U}_t^{(1)}, \mathbf{U}_t^{(2)}, \mathbf{U}_t^{(3)}, \mathbf{G}_t)$ be the t -th iterate of Algorithm 1. Suppose that the induction hypothesis (19) holds at the t -th iteration. Under the assumption that $\alpha \leq \frac{c_1}{\mu^2 r_1 r_2 r_3}$ for some sufficiently small constant c_1 and the choice of ζ_{t+1} in Theorem 1, the $(t+1)$ -th iterate \mathbf{F}_{t+1} satisfies*

$$\max_k \left\{ \sqrt{\frac{n_k}{r_k}} \left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_{t+1}^{(k)} - \mathbf{U}_\star^{(k)}) \Sigma_\star^{(k)} \right\|_{2,\infty} \right\} \leq \rho^{t+1} \sqrt{\mu} \sigma_{\min}(\mathcal{X}_\star)$$

as long as $1/7 \leq \eta \leq 1/4$, where $\{\mathbf{Q}_t^{(k)}\}_{k=1}^3$ are the optimal alignment matrices between \mathbf{F}_t and \mathbf{F}_\star .

It is worthwhile to note that, the local linear convergence of the proposed Algorithm 1, as ensured by the above two lemmas collectively, require the corruption level to not exceed the order of $\frac{1}{\mu^2 r_1 r_2 r_3}$, which is also independent of the condition number. Indeed, the range of the corruption level is mainly constrained by the spectral initialization, as demonstrated next.

Base case: spectral initialization. To establish the induction hypothesis, we still need to check the spectral initialization. The following lemmas state that the spectral initialization satisfies the induction hypothesis (19) at the base case $t = 0$, allowing us to invoke local contraction.

Lemma 3 (Distance at initialization). *Let $\mathbf{Y} = \mathcal{X}_\star + \mathcal{S}_\star \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, where \mathcal{X}_\star is μ -incoherent with rank $\mathbf{r} = (r_1, r_2, r_3)$, and \mathcal{S}_\star is α -sparse. Let $\mathbf{F}_0 := (\mathbf{U}_0^{(1)}, \mathbf{U}_0^{(2)}, \mathbf{U}_0^{(3)}, \mathbf{G}_0)$ be the output of spectral initialization with the threshold obeying $\|\mathcal{X}_\star\|_\infty \leq \zeta_0 \leq 2 \|\mathcal{X}_\star\|_\infty$. If $\alpha \leq \frac{c_0}{\sqrt{\mu^3 r_1 r_2 r_3 r \kappa}}$ for some constant $c_0 > 0$, we have*

$$\text{dist}(\mathbf{F}_0, \mathbf{F}_\star) \leq 54.1 c_0 \sigma_{\min}(\mathcal{X}_\star).$$

Lastly, the next lemma ensures that our initialization satisfies the incoherence condition, which requires nontrivial efforts to exploit the algebraic structures of the Tucker decomposition.

Lemma 4 (Incoherence at initialization). *Let $\mathbf{Y} = \mathcal{X}_\star + \mathcal{S}_\star \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, where \mathcal{X}_\star is μ -incoherent with rank $\mathbf{r} = (r_1, r_2, r_3)$, and \mathcal{S}_\star is α -sparse. Let $\mathbf{F}_0 := (\mathbf{U}_0^{(1)}, \mathbf{U}_0^{(2)}, \mathbf{U}_0^{(3)}, \mathbf{G}_0)$ be the output of spectral*

initialization with the threshold obeying $\|\mathbf{X}_*\|_\infty \leq \zeta_0 \leq 2\|\mathbf{X}_*\|_\infty$. If $\alpha \leq \frac{c_0}{\mu^2 r_1 r_2 r_3 \kappa}$ for some sufficiently small constant c_0 , then the spectral initialization satisfies the incoherence condition

$$\max_k \left\{ \sqrt{\frac{n_k}{r_k}} \left\| (U_0^{(k)} Q_0^{(k)} - U_*^{(k)}) \Sigma_*^{(k)} \right\|_{2,\infty} \right\} \leq \sqrt{\mu} \sigma_{\min}(\mathbf{X}_*),$$

where $\{Q_0^{(k)}\}_{k=1}^3$ are the optimal alignment matrices between \mathbf{F}_0 and \mathbf{F}_* .

4 Numerical experiments

4.1 Experiments on synthetic data

We begin with evaluating the phase transition performance of ScaledGD (cf. Algorithm 1) with respect to the multilinear rank and the level of corruption. For each κ , we randomly generate an $n \times n \times n$ tensor, with $n = 100$, multilinear rank $\mathbf{r} = (r, r, r)$, $r \in \{2, 5, 10, 20, \dots, 80\}$, and level of corruption $\alpha \in \{0.1, 0.2, \dots, 0.9, 1\}$. The factor matrices are generated uniformly at random with orthonormal columns, and a diagonal core tensor \mathcal{G}_* is generated such that $[\mathcal{G}_*]_{i,i,i} = \kappa^{-(i-1)/(r-1)}$ for $i = 1, 2, \dots, r$. We further randomly corrupt α -fraction of the entries, by adding uniformly sampled numbers from the range $[-\sum_{i,j,k} |[\mathbf{X}_*]_{i,j,k}|/n^3, \sum_{i,j,k} |[\mathbf{X}_*]_{i,j,k}|/n^3]$ to the selected entries, where $\sum_{i,j,k} |[\mathbf{X}_*]_{i,j,k}|/n^3$ is the mean of the entry-wise magnitudes of \mathbf{X}_* . To tune the constant step size η , and the hyperparameters ζ_0, ζ_1 , and the decay rate ρ of the thresholding parameter for each tensor automatically, we used the Bayesian optimization method described in [ASY⁺19]. Specifically, we run the toolbox [ASY⁺19] for 200 trials or until the tuned parameters satisfy $\frac{\|\mathbf{X}_T - \mathbf{X}_*\|_F}{\|\mathbf{X}_*\|_F} < 10^{-6}$ for $T = 200$, whichever happened first.

Figure 1 shows the log median of the relative reconstruction error $\frac{\|\mathbf{X}_T - \mathbf{X}_*\|_F}{\|\mathbf{X}_*\|_F}$ when $T = 200$, over 20 random tensor realizations for $\kappa = 1, 5, 10$. Our results show a distinct negative linear relationship between the corruption level and the multilinear rank with regards to the final relative loss of ScaledGD. In particular, the performance is almost independent of the condition number κ , suggesting the performance of ScaledGD is indeed quite insensitive to the condition number.

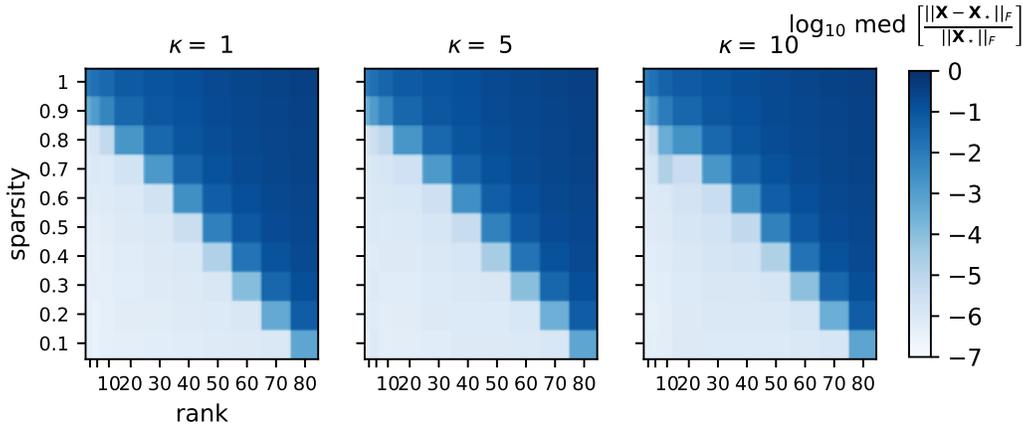


Figure 1: Log median of the relative reconstruction error of $\frac{\|\mathbf{X}_T - \mathbf{X}_*\|_F}{\|\mathbf{X}_*\|_F}$ across 20 randomly generated tensors with varying ranks and levels of corruption when the condition number is set as $\kappa = 1, 5, 10$.

We further investigate the effect of the decay rate ρ of the thresholding parameters while fixing the other hyperparameters tuned as earlier. Using the same method, we generate a $100 \times 100 \times 100$ tensor with $\kappa = 5$ and 20% of the entries corrupted. Figure 2 shows the relative reconstruction error versus the iteration count using different decay rates ρ . It can be seen that ScaledGD enables exact recovery over a wide range of ρ as long as it is not too small. Moreover, within the range of decay rates that still admits exact recovery, the smaller ρ is, the faster ScaledGD converges. Note that the tuned decay rate $\rho \approx 0.931$ does not achieve the

fastest convergence rate since the stopping criteria for hyperparameter tuning were not set to optimize the convergence rate but some accuracy-speed trade-off.

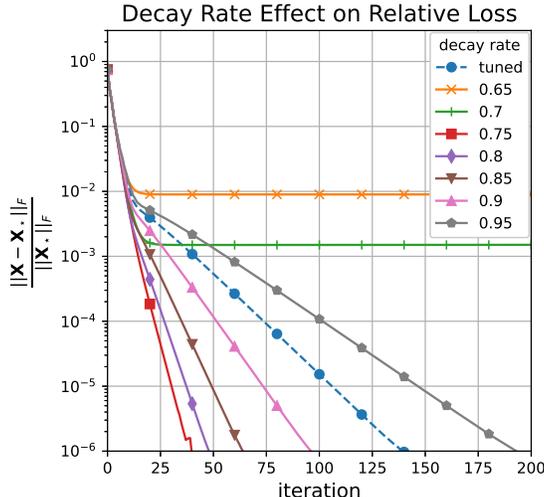


Figure 2: The relative reconstruction error $\frac{\|\mathcal{X}_T - \mathcal{X}_*\|_F}{\|\mathcal{X}_*\|_F}$ with respect to the iteration count, when varying the decay rate ρ with other hyperparameters fixed.

Next, we also examine the performance of ScaledGD with shot noise, corruptions drawn from a Poisson distribution, with comparisons to the Riemannian gradient descent (RiemannianGD) algorithm in [CLX21]. More specifically, if corrupted, a ground truth entry $\mathcal{X}_{i,j,k}$ has noise drawn from $10^{-5}\text{Poisson}(10^5|\mathcal{X}_{i,j,k}|)$ added to itself, where the 10^5 scaling is to encourage draws that are nonzero. This type of noise is non-negative and perturbs higher magnitude entries more. Note that the per-iteration cost of RiemannianGD is significantly higher than ours, due to the fact that it requires an evaluation of a rank- \mathbf{r} HOSVD. Therefore, for this experiment, we generate tensors of a smaller size $50 \times 50 \times 50$ to accommodate the high computation need of RiemannianGD. We similarly tune the hyperparameters of RiemannianGD using the same method mentioned earlier for 100 trials. Figure 3 shows the log median of the relative reconstruction error $\frac{\|\mathcal{X}_T - \mathcal{X}_*\|_F}{\|\mathcal{X}_*\|_F}$ when $T = 100$, over 20 random tensor realizations when $\kappa = 5$. It can be observed that the empirical performance of the two methods, indicated by the phase transition curves, are similar when tuned properly. However, the ScaledGD method is considerably faster, the difference of which is accentuated for larger and lower rank tensors, due to the fact that it works in the factor space and does not need to perform rank- \mathbf{r} HOSVD at every iteration.

4.2 Image denoising and outlier detection

In this experiment, we examine the performance of ScaledGD for imaging denoising and outlier detection, with comparisons to the tensor RPCA algorithm proposed in [LFC+19] called TNN for their use of a newly defined tensor nuclear norm (TNN). We consider a sequence of handwritten digits “2” from the MNIST database [LCB10] containing 5958 images of size 28×28 , leading to a 3-way tensor. We assume the tensor is low-rank along the image sequence, but not within the image for simplicity; in other words, the multilinear rank is assumed as $\mathbf{r} = (5, 28, 28)$. For both algorithms, the hyperparameters are best tuned by hands.

We examine the performance of ScaledGD and TNN when the image sequence is contaminated in the following scenarios: 1) 70% salt and pepper noise; 2) 500 out of the total images are randomly selected and swapped by random images from the entire MNIST training set; and 3) 50% salt and pepper noise and 500 randomly swapped images. Figure 4 demonstrates the performance of the compared algorithms on the first 100 instances for each corruption scenario. In all situations, ScaledGD recovers the low-rank component corresponding to the correct digit more accurately than TNN from a visual inspection. Furthermore, ScaledGD corrected the oddly-shaped or outlying digits to make the low-rank component be more homogeneous, but

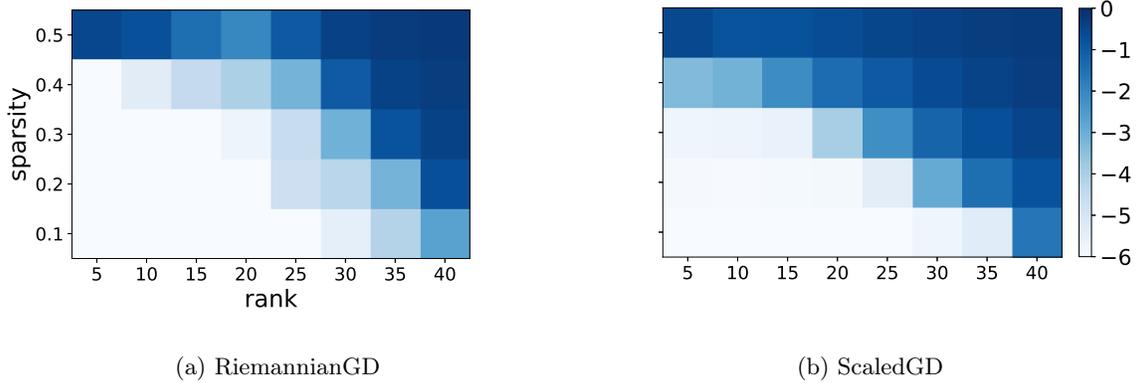


Figure 3: Comparison between RiemannianGD [CLX21] and ScaledGD on the log median of the relative recovery error, $\frac{\|\mathbf{x}^* - \mathbf{x}_T\|_F}{\|\mathbf{x}^*\|_F}$, across 20 randomly generated tensors with varying ranks and levels of shot noise corruption when the condition number is set as $\kappa = 5$.

TNN mostly preserved these cases in the low-rank output. More importantly, ScaledGD runs much faster as a scalable nonconvex approach, while TNN is more computationally expensive using convex optimization.

4.3 Background subtraction via selective mode update

We now apply ScaledGD to the task of background subtraction using videos from the VIRAT dataset [OHP⁺11], where the height and width of the videos are reduced by a factor of 4 due to hardware limitations. The video data can be thought as a multi-way tensor spanning across the height, width, frames, as well as different color channels of the scene. Here, the low-rank tensor corresponds to the background in the video which is fairly static over the frames, and the sparse tensor corresponds to the foreground containing moving objects which takes a small number of pixels. In particular, it is reasonable to assume that the background tensor is low-rank for the mode corresponding to the frames, but full rank in other modes. Motivated by this observation, one might be tempted to selectively only update the core tensor and the factor matrix corresponding to the mode for frames while keeping the other factor matrices fixed after the spectral initialization. We compare the results using this selective mode update strategy with the original ScaledGD algorithm in Figure 5, where the same hyperparameters are used in both. It can be seen that skipping updates of the full-rank factor matrices produced qualitatively similar results while gaining a significant per-iteration speed-up of about 4.6 to 5 times. We expect the speed improvement to be greater for larger tensors, as more computation can be bypassed.

5 Conclusions

In this paper, we proposed a new algorithm for tensor RPCA, which applies scaled gradient descent with an iteration-varying thresholding operation to adaptively remove the impact of corruptions. The algorithm is demonstrated to be fast, provably accurate, and achieve competitive performance over real-world applications. It opens several interesting directions for future research.

- *Dependency with the condition number from spectral initialization.* As seen from the analysis, the local linear convergence of Algorithm 1 succeeds under a larger range of the sparsity level α independent of the condition number κ . The constraint on α with respect to the condition number κ mainly stems from the spectral initialization, and it is of great interest to see if it is possible to refine the analysis in terms of the dependency with κ , which likely will require new tools.

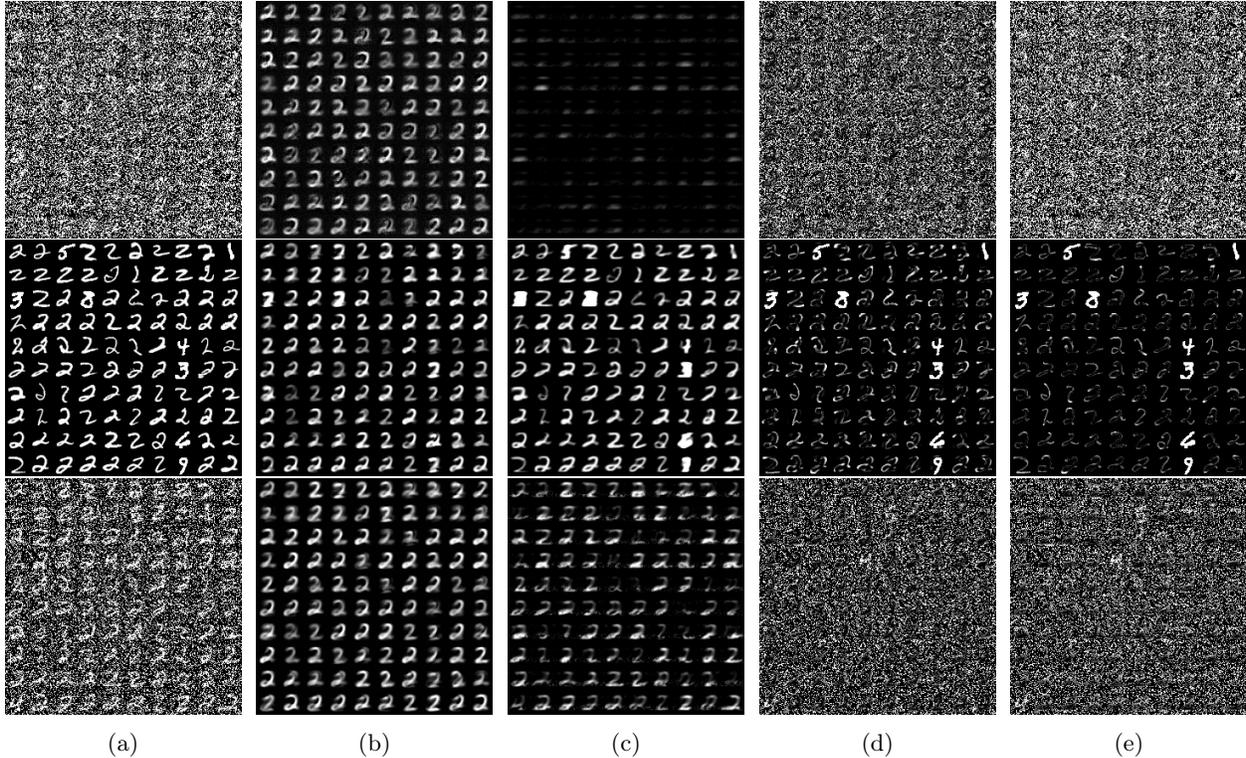


Figure 4: Imaging denoising and outlier removal on an image sequence of handwritten digits with various corruption scenarios, using tensor RPCA via ScaledGD and TNN [LFC⁺19]. From top to bottom, the rows show results on the first 100 images when 1) corrupted with 70% salt and pepper noise; 2) 500 randomly swapped images; and 3) 50% salt and pepper noise and 500 randomly swapped images. From left to right, the columns show (a) the corrupted input, (b) the low-rank output of ScaledGD, (c) the low-rank output of TNN, (d) the sparse output of ScaledGD, and (e) the sparse output of TNN.

- *Missing data.* An important extension is to handle missing data in tensor RPCA, which seeks to recover a low-rank ground truth tensor from its partially observed entries possibly corrupted by gross errors. Our proposed algorithm can be adapted to this case in a straightforward fashion by considering the loss function defined only over the observed entries, and understanding its performance guarantees is a natural step.
- *Streaming data.* An equally interesting direction is to perform tensor RPCA over online and streaming data, where the fibers or slices of the tensor arrive sequentially over time, a situation that is common in data analytics [BCL18, VN18]. It is of great interest to develop low-complexity algorithms that can estimate and track the low-rank tensor factors as quickly as possible.
- *Hyperparameters.* The proposed algorithm contains several hyperparameters that need to be tuned carefully to fully unleash its potential. A recent follow-up [DSDC22] examined a learned approach based on algorithm unfolding and self-supervised learning to enable automatic hyperparameter tuning. In addition, understanding the performance when the rank is only imperfectly specified is also of great importance, which is closely related to [XSCM23].

Acknowledgements

The work of H. Dong, T. Tong and Y. Chi is supported in part by Office of Naval Research under N00014-19-1-2404, by Air Force Research Laboratory under FA8750-20-2-0504, by National Science Foundation under CAREER ECCS-1818571, CCF-1901199 and ECCS-2126634, and by Department of Transportation under

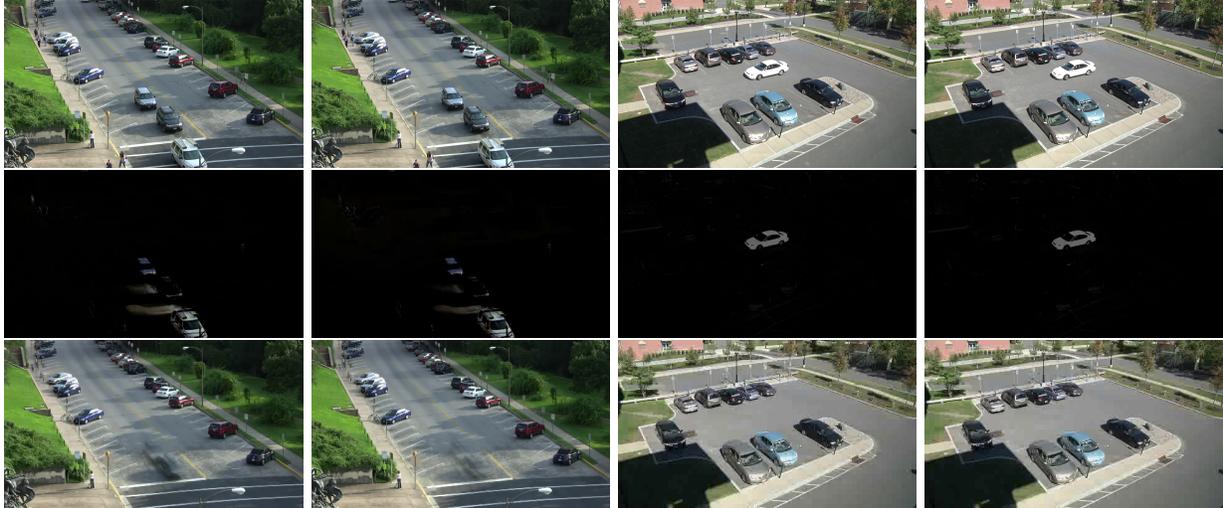


Figure 5: Examples of extracted background and foreground in video surveillance with and without selective mode updates. The first two columns use the same street surveillance video, and the last two columns use the same parking lot video, where the same frame is shown for each pair. In a pair, the left column is the result when using the original ScaledGD algorithm, and the right column is the result employing selective mode updates. The first row is the original frame, and second row is the sparse foreground, and the third row is the background.

693JJ321C000013. The work of H. Dong is also supported by CIT Dean’s Fellowship, Liang Ji-Dian Graduate Fellowship, and Michel and Kathy Doreau Graduate Fellowship in Electrical and Computer Engineering at Carnegie Mellon University. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon.

References

- [AGH⁺14] A. Anandkumar, R. Ge, D. Hsu, S. Kakade, and M. Telgarsky. Tensor decompositions for learning latent variable models. *Journal of Machine Learning Research*, 15:2773–2832, 2014.
- [AJSN16] A. Anandkumar, P. Jain, Y. Shi, and U. N. Niranjan. Tensor vs. matrix methods: Robust tensor decomposition under block sparse perturbations. In *Artificial Intelligence and Statistics*, pages 268–276. PMLR, 2016.
- [ARB20] T. Ahmed, H. Raja, and W. U. Bajwa. Tensor regression using low-rank and sparse Tucker decompositions. *SIAM Journal on Mathematics of Data Science*, 2(4):944–966, 2020.
- [ASY⁺19] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 2623–2631, 2019.
- [BCL18] L. Balzano, Y. Chi, and Y. M. Lu. Streaming PCA and subspace tracking: The missing data case. *Proceedings of the IEEE*, 106(8):1293–1310, 2018.
- [BL10] G. Bergqvist and E. G. Larsson. The higher-order singular value decomposition: Theory and an application [lecture notes]. *IEEE Signal Processing Magazine*, 27(3):151–154, 2010.
- [CC14] Y. Chen and Y. Chi. Robust spectral compressed sensing via structured matrix completion. *IEEE Transactions on Information Theory*, 60(10):6576 – 6601, Oct. 2014.

- [CC18] Y. Chen and Y. Chi. Harnessing structures in big data via guaranteed low-rank matrix estimation: Recent theory and fast algorithms via convex and nonconvex optimization. *IEEE Signal Processing Magazine*, 35(4):14 – 31, 2018.
- [CCW19] H. Cai, J.-F. Cai, and K. Wei. Accelerated alternating projections for robust principal component analysis. *The Journal of Machine Learning Research*, 20(1):685–717, 2019.
- [CFMY21] Y. Chen, J. Fan, C. Ma, and Y. Yan. Bridging convex and nonconvex optimization in robust PCA: Noise, outliers, and missing data. *The Annals of Statistics*, 49(5):2948–2971, 2021.
- [CLC19] Y. Chi, Y. M. Lu, and Y. Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing*, 67(20):5239–5269, 2019.
- [CLMW11] E. J. Candès, X. Li, Y. Ma, and J. Wright. Robust principal component analysis? *Journal of the ACM*, 58(3):11:1–11:37, 2011.
- [CLX21] J.-F. Cai, J. Li, and D. Xia. Generalized low-rank plus sparse tensor estimation by fast Riemannian optimization. *arXiv preprint arXiv:2103.08895*, 2021.
- [CLY21] H. Cai, J. Liu, and W. Yin. Learned robust PCA: A scalable deep unfolding approach for high-dimensional outlier detection. *Advances in Neural Information Processing Systems*, 34, 2021.
- [CSPW11] V. Chandrasekaran, S. Sanghavi, P. Parrilo, and A. Willsky. Rank-sparsity incoherence for matrix decomposition. *SIAM Journal on Optimization*, 21(2):572–596, 2011.
- [DBBG19] D. Driggs, S. Becker, and J. Boyd-Graber. Tensor robust principal component analysis: Better recovery with atomic norm regularization. *arXiv preprint arXiv:1901.10991*, 2019.
- [DLDMV00] L. De Lathauwer, B. De Moor, and J. Vandewalle. A multilinear singular value decomposition. *SIAM journal on Matrix Analysis and Applications*, 21(4):1253–1278, 2000.
- [DSDC22] H. Dong, M. Shah, S. Donegan, and Y. Chi. Deep unfolded tensor robust PCA with self-supervised learning. *arXiv preprint arXiv:2212.11346*, 2022.
- [FL18] S. Friedland and L.-H. Lim. Nuclear norm of higher-order tensors. *Mathematics of Computation*, 87(311):1255–1281, 2018.
- [GGH14] Q. Gu, H. Gui, and J. Han. Robust tensor decomposition with gross corruption. *Advances in Neural Information Processing Systems*, 27, 2014.
- [GQ14] D. Goldfarb and Z. Qin. Robust low-rank tensor recovery: Models and algorithms. *SIAM Journal on Matrix Analysis and Applications*, 35(1):225–253, 2014.
- [GWL16] Q. Gu, Z. Wang, and H. Liu. Low-rank and sparse structure pursuit via alternating minimization. In *Artificial Intelligence and Statistics*, pages 600–609, 2016.
- [HWZ20] R. Han, R. Willett, and A. Zhang. An optimal statistical and computational framework for generalized tensor estimation. *arXiv preprint arXiv:2002.11255*, 2020.
- [KABO10] A. Karatzoglou, X. Amatriain, L. Baltrunas, and N. Oliver. Multiverse recommendation: n -dimensional tensor factorization for context-aware collaborative filtering. In *Proceedings of the fourth ACM conference on Recommender systems*, pages 79–86, 2010.
- [KB09] T. G. Kolda and B. W. Bader. Tensor decompositions and applications. *SIAM review*, 51(3):455–500, 2009.
- [LCB10] Y. LeCun, C. Cortes, and C. Burges. Mnist handwritten digit database. *ATT Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.

- [LCM10] Z. Lin, M. Chen, and Y. Ma. The augmented Lagrange multiplier method for exact recovery of corrupted low-rank matrices. *arXiv preprint arXiv:1009.5055*, 2010.
- [LFC⁺16] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5249–5257, 2016.
- [LFC⁺19] C. Lu, J. Feng, Y. Chen, W. Liu, Z. Lin, and S. Yan. Tensor robust principal component analysis with a new tensor nuclear norm. *IEEE transactions on pattern analysis and machine intelligence*, 42(4):925–938, 2019.
- [LMWY12] J. Liu, P. Musialski, P. Wonka, and J. Ye. Tensor completion for estimating missing values in visual data. *IEEE transactions on Pattern Analysis and Machine Intelligence*, 35(1):208–220, 2012.
- [LWQ⁺15] S. Li, W. Wang, H. Qi, B. Ayhan, C. Kwan, and S. Vance. Low-rank tensor decomposition based anomaly detection for hyperspectral imagery. In *2015 IEEE International Conference on Image Processing (ICIP)*, pages 4525–4529. IEEE, 2015.
- [LZ21] Y. Luo and A. R. Zhang. Low-rank tensor estimation via Riemannian Gauss-Newton: Statistical optimality and second-order convergence. *arXiv preprint arXiv:2104.12031*, 2021.
- [NNS⁺14] P. Netrapalli, U. Niranjan, S. Sanghavi, A. Anandkumar, and P. Jain. Non-convex robust PCA. In *Advances in Neural Information Processing Systems*, pages 1107–1115, 2014.
- [OHP⁺11] S. Oh, A. Hoogs, A. Perera, N. Cuntoor, C.-C. Chen, J. T. Lee, S. Mukherjee, J. Aggarwal, H. Lee, L. Davis, et al. A large-scale benchmark dataset for event recognition in surveillance video. In *Proceedings of the 2011 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3153–3160, 2011.
- [PFS16] E. E. Papalexakis, C. Faloutsos, and N. D. Sidiropoulos. Tensors for data mining and data fusion: Models, applications, and scalable algorithms. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 8(2):1–44, 2016.
- [TMC21a] T. Tong, C. Ma, and Y. Chi. Accelerating ill-conditioned low-rank matrix estimation via scaled gradient descent. *Journal of Machine Learning Research*, 22(150):1–63, 2021.
- [TMC21b] T. Tong, C. Ma, and Y. Chi. Low-rank matrix recovery with scaled subgradient methods: Fast and robust convergence without the condition number. *IEEE Transactions on Signal Processing*, 69:2396–2409, 2021.
- [TMC22] T. Tong, C. Ma, and Y. Chi. Accelerating ill-conditioned robust low-rank tensor regression. In *International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*. IEEE, 2022.
- [TMPB⁺22] T. Tong, C. Ma, A. Prater-Bennette, E. Tripp, and Y. Chi. Scaling and scalability: Provable nonconvex low-rank tensor estimation from incomplete measurements. *Journal of Machine Learning Research*, 23(163):1–77, 2022.
- [TY10] K.-C. Toh and S. Yun. An accelerated proximal gradient algorithm for nuclear norm regularized linear least squares problems. *Pacific Journal of optimization*, 6(615-640):15, 2010.
- [VN18] N. Vaswani and P. Narayanamurthy. Static and dynamic robust PCA and matrix completion: A review. *Proceedings of the IEEE*, 106(8):1359–1379, 2018.
- [VVM12] N. Vannieuwenhoven, R. Vandebril, and K. Meerbergen. A new truncation strategy for the higher-order singular value decomposition. *SIAM Journal on Scientific Computing*, 34(2):A1027–A1052, 2012.
- [WCW21] H. Wang, J. Chen, and K. Wei. Entrywise convergence of riemannian gradient method for low rank tensor completion via tucker decomposition. *arXiv preprint arXiv:2108.07899*, 2021.

- [WGR⁺09] J. Wright, A. Ganesh, S. Rao, Y. Peng, and Y. Ma. Robust principal component analysis: Exact recovery of corrupted low-rank matrices via convex optimization. *Advances in neural information processing systems*, 22, 2009.
- [XSCM23] X. Xu, Y. Shen, Y. Chi, and C. Ma. The power of preconditioning in overparameterized low-rank matrix sensing. *arXiv preprint arXiv:2302.01186*, 2023.
- [XY19] D. Xia and M. Yuan. On polynomial time methods for exact low-rank tensor completion. *Foundations of Computational Mathematics*, 19(6):1265–1313, 2019.
- [YPCC16] X. Yi, D. Park, Y. Chen, and C. Caramanis. Fast algorithms for robust PCA via gradient descent. In *Advances in Neural Information Processing Systems*, pages 4152–4160, 2016.
- [YZ16] M. Yuan and C.-H. Zhang. On tensor completion via nuclear norm minimization. *Foundations of Computational Mathematics*, 16(4):1031–1068, 2016.
- [ZWZM19] X. Zhang, D. Wang, Z. Zhou, and Y. Ma. Robust low-rank tensor recovery with rectification and alignment. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(1):238–255, 2019.

A Proof of Lemma 1

Since $\text{dist}(\mathbf{F}_t, \mathbf{F}_\star) < \sigma_{\min}(\mathcal{X}_\star)$, [TMPB⁺22, Lemma 6] ensures that the optimal alignment matrices $\{\mathbf{Q}_t^{(k)}\}_{k=1}^3$ between \mathbf{F}_t and \mathbf{F}_\star exist. Since $\{\mathbf{Q}_t^{(k)}\}_{k=1}^3$ may be a suboptimal alignment between \mathbf{F}_{t+1} and \mathbf{F}_\star , we therefore have

$$\text{dist}^2(\mathbf{F}_{t+1}, \mathbf{F}_\star) \leq \sum_{k=1}^3 \left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\|_{\text{F}}^2 + \left\| ((\mathbf{Q}_t^{(1)})^{-1}, (\mathbf{Q}_t^{(2)})^{-1}, (\mathbf{Q}_t^{(3)})^{-1}) \cdot \mathbf{g}_{t+1} - \mathbf{g}_\star \right\|_{\text{F}}^2. \quad (20)$$

Before we embark on the control of the terms on the right hand side of (20), we introduce the following short-hand notations:

$$\begin{aligned} \mathbf{U}^{(k)} &:= \mathbf{U}_t^{(k)} \mathbf{Q}_t^{(k)}, & \check{\mathbf{U}}^{(k)} &:= \check{\mathbf{U}}_t^{(k)} (\mathbf{Q}_t^{(k)})^{-\top}, & \mathbf{S} &:= \mathbf{S}_{t+1}, \\ \Delta_{\mathbf{U}^{(k)}} &:= \mathbf{U}^{(k)} - \mathbf{U}_\star^{(k)}, & \Delta_{\check{\mathbf{U}}^{(k)}} &:= \check{\mathbf{U}}^{(k)} - \check{\mathbf{U}}_\star^{(k)}, & \Delta_{\mathbf{S}} &:= \mathbf{S} - \mathbf{S}_\star, \\ \mathbf{g} &:= ((\mathbf{Q}_t^{(1)})^{-1}, (\mathbf{Q}_t^{(2)})^{-1}, (\mathbf{Q}_t^{(3)})^{-1}) \cdot \mathbf{g}_t, & \Delta_{\mathbf{g}} &:= \mathbf{g} - \mathbf{g}_\star. \end{aligned} \quad (21)$$

In addition, Lemma 10 in conjunction with the induction hypothesis (19) tells us that

$$\|\mathcal{X}_t - \mathcal{X}_\star\|_\infty \leq 8 \sqrt{\frac{\mu^3 r_1 r_2 r_3}{n_1 n_2 n_3}} \rho^t \sigma_{\min}(\mathcal{X}_\star) =: \zeta_{t+1}. \quad (22)$$

Step 1: bounding the first term of (20). Using the update rule (16b) for $\mathbf{U}_{t+1}^{(k)}$, we have

$$\begin{aligned} & (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \\ &= \left[\left((1-\eta) \mathbf{U}_t^{(k)} - \eta (\mathcal{M}_k(\mathbf{S}_{t+1}) - \mathcal{M}_k(\mathbf{y})) \check{\mathbf{U}}_t^{(k)} (\check{\mathbf{U}}_t^{(k)\top} \check{\mathbf{U}}_t^{(k)})^{-1} \right) \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)} \right] \boldsymbol{\Sigma}_\star^{(k)} \\ &= \left[\left((1-\eta) \mathbf{U}_t^{(k)} - \eta (\mathcal{M}_k(\mathbf{S}_{t+1} - \mathbf{S}_\star) - \mathbf{U}_\star^{(k)} \check{\mathbf{U}}_\star^{(k)\top}) \check{\mathbf{U}}_t^{(k)} (\check{\mathbf{U}}_t^{(k)\top} \check{\mathbf{U}}_t^{(k)})^{-1} \right) \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)} \right] \boldsymbol{\Sigma}_\star^{(k)}, \end{aligned}$$

where the second equality follows from $\mathbf{y} = \mathcal{X}_\star + \mathbf{S}_\star$ as well as the matricization property (4). With the set of notation (21) in place, simple algebraic simplifications yield

$$\begin{aligned} (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} &= \left[(1-\eta) \mathbf{U}^{(k)} - \eta (\mathcal{M}_k(\Delta_{\mathbf{S}}) - \mathbf{U}_\star^{(k)} \check{\mathbf{U}}_\star^{(k)\top}) \check{\mathbf{U}}^{(k)} (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} - \mathbf{U}_\star^{(k)} \right] \boldsymbol{\Sigma}_\star^{(k)} \\ &= \left[(1-\eta) \Delta_{\mathbf{U}^{(k)}} - \eta (\mathcal{M}_k(\Delta_{\mathbf{S}}) + \mathbf{U}_\star^{(k)} \Delta_{\check{\mathbf{U}}^{(k)}}^\top) \check{\mathbf{U}}^{(k)} (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} \right] \boldsymbol{\Sigma}_\star^{(k)} \end{aligned}$$

$$= (1 - \eta) \mathbf{\Delta}_{U^{(k)}} \mathbf{\Sigma}_*^{(k)} - \eta \left(\mathcal{M}_k(\mathbf{\Delta}_S) + \mathbf{U}_*^{(k)} \mathbf{\Delta}_{U^{(k)}}^\top \right) \check{\mathbf{U}}^{(k)} (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} \mathbf{\Sigma}_*^{(k)}. \quad (23)$$

Detailed in Appendix A.1, we claim the following bound holds:

$$\begin{aligned} \sum_{k=1}^3 \left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_*^{(k)}) \mathbf{\Sigma}_*^{(k)} \right\|_{\mathbb{F}}^2 &\leq (1 - \eta)^2 \sum_{k=1}^3 \left\| \mathbf{\Delta}_{U^{(k)}} \mathbf{\Sigma}_*^{(k)} \right\|_{\mathbb{F}}^2 + 0.15\eta(1 - \eta) \sum_{k=1}^3 \left\| \mathbf{\Delta}_{U^{(k)}} \mathbf{\Sigma}_*^{(k)} \right\|_{\mathbb{F}} \epsilon_0 \rho^t \sigma_{\min}(\mathbf{\mathcal{X}}_*) \\ &\quad + 2\eta^2 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathbf{\mathcal{X}}_*) + 0.06\eta(1 - \eta) \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathbf{\mathcal{X}}_*). \end{aligned} \quad (24)$$

Step 2: bounding the second term of (20). Using the update rule (16c) for \mathbf{g}_{t+1} , we have

$$\begin{aligned} &((\mathbf{Q}_t^{(1)})^{-1}, (\mathbf{Q}_t^{(2)})^{-1}, (\mathbf{Q}_t^{(3)})^{-1}) \cdot \mathbf{g}_{t+1} - \mathbf{g}_* \\ &= (1 - \eta) ((\mathbf{Q}_t^{(1)})^{-1}, (\mathbf{Q}_t^{(2)})^{-1}, (\mathbf{Q}_t^{(3)})^{-1}) \cdot \mathbf{g}_t \\ &\quad - \eta \left((\mathbf{U}_t^{(1)\top} \mathbf{U}_t^{(1)} \mathbf{Q}_t^{(1)})^{-1} \mathbf{U}_t^{(1)\top}, (\mathbf{U}_t^{(2)\top} \mathbf{U}_t^{(2)} \mathbf{Q}_t^{(2)})^{-1} \mathbf{U}_t^{(2)\top}, (\mathbf{U}_t^{(3)\top} \mathbf{U}_t^{(3)} \mathbf{Q}_t^{(3)})^{-1} \mathbf{U}_t^{(3)\top} \right) \cdot (\mathbf{s}_{t+1} - \mathbf{y}) - \mathbf{g}_* \\ &= (1 - \eta) \mathbf{g} - \eta \left((\mathbf{U}^{(1)\top} \mathbf{U}^{(1)})^{-1} \mathbf{U}^{(1)\top}, (\mathbf{U}^{(2)\top} \mathbf{U}^{(2)})^{-1} \mathbf{U}^{(2)\top}, (\mathbf{U}^{(3)\top} \mathbf{U}^{(3)})^{-1} \mathbf{U}^{(3)\top} \right) \cdot (\mathbf{\Delta}_S - \mathbf{\mathcal{X}}_*) - \mathbf{g}_* \\ &= (1 - \eta) \mathbf{\Delta}_g - \eta \left((\mathbf{U}^{(1)\top} \mathbf{U}^{(1)})^{-1} \mathbf{U}^{(1)\top}, (\mathbf{U}^{(2)\top} \mathbf{U}^{(2)})^{-1} \mathbf{U}^{(2)\top}, (\mathbf{U}^{(3)\top} \mathbf{U}^{(3)})^{-1} \mathbf{U}^{(3)\top} \right) \\ &\quad \cdot \left((\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathbf{g}_* - \mathbf{\mathcal{X}}_* + \mathbf{\Delta}_S \right), \end{aligned} \quad (25)$$

where the last two lines make use of the short-hand notation in (21), as well as the multilinearity property (3). Detailed in Appendix A.2, we claim the following bound holds:

$$\begin{aligned} \left\| ((\mathbf{Q}_t^{(1)})^{-1}, (\mathbf{Q}_t^{(2)})^{-1}, (\mathbf{Q}_t^{(3)})^{-1}) \cdot \mathbf{g}_{t+1} - \mathbf{g}_* \right\|_{\mathbb{F}}^2 &\leq (1 - \eta)^2 \|\mathbf{\Delta}_g\|_{\mathbb{F}}^2 + 2 \cdot 0.15\eta(1 - \eta) \|\mathbf{\Delta}_g\|_{\mathbb{F}} \epsilon_0 \rho^t \sigma_{\min}(\mathbf{\mathcal{X}}_*) \\ &\quad + 0.02\eta(1 - \eta) \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathbf{\mathcal{X}}_*) + 0.06\eta^2 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathbf{\mathcal{X}}_*). \end{aligned} \quad (26)$$

Step 3: putting the bounds together. Plugging the bounds (24) and (26) into (20) yields

$$\begin{aligned} &\text{dist}^2(\mathbf{F}_{t+1}, \mathbf{F}_*) \\ &\leq (1 - \eta)^2 \left(\sum_{k=1}^3 \left\| \mathbf{\Delta}_{U^{(k)}} \mathbf{\Sigma}_*^{(k)} \right\|_{\mathbb{F}}^2 + \|\mathbf{\Delta}_g\|_{\mathbb{F}}^2 \right) + 0.15\eta(1 - \eta) \left(\sum_{k=1}^3 \left\| \mathbf{\Delta}_{U^{(k)}} \mathbf{\Sigma}_*^{(k)} \right\|_{\mathbb{F}} + 2 \|\mathbf{\Delta}_g\|_{\mathbb{F}} \right) \epsilon_0 \rho^t \sigma_{\min}(\mathbf{\mathcal{X}}_*) \\ &\quad + 2.1\eta^2 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathbf{\mathcal{X}}_*) + 0.08\eta(1 - \eta) \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathbf{\mathcal{X}}_*) \\ &\stackrel{(i)}{\leq} (1 - \eta)^2 \text{dist}^2(\mathbf{F}_t, \mathbf{F}_*) + 0.4\eta(1 - \eta) \text{dist}(\mathbf{F}_t, \mathbf{F}_*) \epsilon_0 \rho^t \sigma_{\min}(\mathbf{\mathcal{X}}_*) + 2.1\eta^2 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathbf{\mathcal{X}}_*) \\ &\quad + 0.08\eta(1 - \eta) \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathbf{\mathcal{X}}_*) \end{aligned} \quad (27)$$

$$\begin{aligned} &\stackrel{(ii)}{\leq} ((1 - \eta)^2 + 0.4\eta(1 - \eta) + 2.1\eta^2 + 0.08\eta(1 - \eta)) \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathbf{\mathcal{X}}_*) \\ &= ((1 - \eta)^2 + 0.5\eta(1 - \eta) + 2.1\eta^2) \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathbf{\mathcal{X}}_*), \end{aligned} \quad (28)$$

where (i) follows from the definition of $\text{dist}^2(\mathbf{F}_t, \mathbf{F}_*)$ and Cauchy-Schwarz, and (ii) follows from the induction hypothesis $\text{dist}(\mathbf{F}_t, \mathbf{F}_*) \leq \epsilon_0 \rho^t \sigma_{\min}(\mathbf{\mathcal{X}}_*)$. For $0 < \eta \leq 1/4$ and $\rho = 1 - 0.45\eta$, this simplifies to the claimed bound

$$\text{dist}^2(\mathbf{F}_{t+1}, \mathbf{F}_*) \leq (1 - 0.45\eta)^2 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathbf{\mathcal{X}}_*) = \epsilon_0^2 \rho^{2t+2} \sigma_{\min}^2(\mathbf{\mathcal{X}}_*).$$

A.1 Proof of (24)

Taking the squared norm on both sides of (23), we obtain

$$\left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_{t,k} - \mathbf{U}_*^{(k)}) \mathbf{\Sigma}_*^{(k)} \right\|_{\mathbb{F}}^2 = (1 - \eta)^2 \left\| \mathbf{\Delta}_{U^{(k)}} \mathbf{\Sigma}_*^{(k)} \right\|_{\mathbb{F}}^2 + \eta^2 \underbrace{\left\| \mathcal{M}_k(\mathbf{\Delta}_S) \check{\mathbf{U}}^{(k)} (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} \mathbf{\Sigma}_*^{(k)} \right\|_{\mathbb{F}}^2}_{=:\mathfrak{A}_{1,k}}$$

$$\begin{aligned}
& - 2\eta(1-\eta) \underbrace{\left\langle \Delta_{U^{(k)}} \Sigma_\star^{(k)}, U_\star^{(k)} \Delta_{\check{U}^{(k)}}^\top \check{U}^{(k)} (\check{U}^{(k)\top} \check{U}^{(k)})^{-1} \Sigma_\star^{(k)} \right\rangle}_{=:\mathfrak{A}_{2,k}} \\
& + \eta^2 \left\| \underbrace{U_\star^{(k)} \Delta_{\check{U}^{(k)}}^\top \check{U}^{(k)} (\check{U}^{(k)\top} \check{U}^{(k)})^{-1} \Sigma_\star^{(k)}}_{=:\mathfrak{A}_{3,k}} \right\|_{\text{F}}^2 \\
& - 2\eta(1-\eta) \underbrace{\left\langle \Delta_{U^{(k)}} \Sigma_\star^{(k)}, \mathcal{M}_k(\Delta_{\mathcal{S}}) \check{U}^{(k)} (\check{U}^{(k)\top} \check{U}^{(k)})^{-1} \Sigma_\star^{(k)} \right\rangle}_{=:\mathfrak{A}_{4,k}} \\
& + 2\eta^2 \underbrace{\left\langle \mathcal{M}_k(\Delta_{\mathcal{S}}) \check{U}^{(k)} (\check{U}^{(k)\top} \check{U}^{(k)})^{-1} \Sigma_\star^{(k)}, U_\star^{(k)} \Delta_{\check{U}^{(k)}}^\top \check{U}^{(k)} (\check{U}^{(k)\top} \check{U}^{(k)})^{-1} \Sigma_\star^{(k)} \right\rangle}_{=:\mathfrak{A}_{5,k}}.
\end{aligned} \tag{29}$$

In the sequel, we shall bound each term separately.

- **Bounding $\mathfrak{A}_{1,k}$.** Since the quantity inside the norm is of rank r_k , we have

$$\begin{aligned}
\mathfrak{A}_{1,k} &= \left\| \mathcal{M}_k(\Delta_{\mathcal{S}}) \check{U}^{(k)} (\check{U}^{(k)\top} \check{U}^{(k)})^{-1} \Sigma_\star^{(k)} \right\|_{\text{F}}^2 \\
&\leq r_k \left\| \mathcal{M}_k(\Delta_{\mathcal{S}}) \check{U}^{(k)} (\check{U}^{(k)\top} \check{U}^{(k)})^{-1} \Sigma_\star^{(k)} \right\|^2 \\
&\leq r_k \|\mathcal{M}_k(\Delta_{\mathcal{S}})\|^2 \left\| \check{U}^{(k)} (\check{U}^{(k)\top} \check{U}^{(k)})^{-1} \Sigma_\star^{(k)} \right\|^2 \leq \frac{r_k}{(1-\epsilon_0)^6} \|\mathcal{M}_k(\Delta_{\mathcal{S}})\|^2,
\end{aligned}$$

where the last inequality follows from Lemma 8 (cf. (70c)). To continue, notice that the choice of ζ_{t+1} (cf. (22)) guarantees that $\Delta_{\mathcal{S}}$ (and hence $\mathcal{M}_k(\Delta_{\mathcal{S}})$) is α -sparse (cf. Lemma 12). This allows us to invoke Lemma 11 and obtain

$$\|\mathcal{M}_k(\Delta_{\mathcal{S}})\| \leq \alpha \sqrt{n_1 n_2 n_3} \|\mathcal{M}_k(\Delta_{\mathcal{S}})\|_\infty = \alpha \sqrt{n_1 n_2 n_3} \|\Delta_{\mathcal{S}}\|_\infty \leq 2\alpha \sqrt{n_1 n_2 n_3} \zeta_{t+1}. \tag{30}$$

Plugging this into the previous inequality, we arrive at

$$\mathfrak{A}_{1,k} \leq \frac{4\alpha^2 n_1 n_2 n_3 r_k}{(1-\epsilon_0)^6} \zeta_{t+1}^2 \leq \frac{256\alpha^2 \mu^3 r_1 r_2 r_3 r_k}{(1-\epsilon_0)^6} \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star),$$

where the second inequality follows from the choice of ζ_{t+1} (cf. (22)). Finally, with the assumption on the sparsity level $\alpha \leq \frac{c_0 \epsilon_0}{\sqrt{\mu^3 r_1 r_2 r_3 r}}$, we have

$$\mathfrak{A}_{1,k} \leq \frac{256c_0^2}{(1-\epsilon_0)^6} \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star) \leq 0.02 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star) \tag{31}$$

for sufficiently small c_0 and $\epsilon_0 < 0.01$.

- **Bounding $\mathfrak{A}_{2,k}$.** This term is identical to the term that is bounded in [TMPB⁺22, Section B.1], which obeys

$$\mathfrak{A}_{2,k} \geq \left\langle \mathcal{K}^{(k)}, \mathcal{K}^{(1)} + \mathcal{K}^{(2)} + \mathcal{K}^{(3)} \right\rangle - C_1 \epsilon_0 \rho^t \text{dist}^2(\mathbf{F}_t, \mathbf{F}_\star)$$

for some constant $C_1 > 1$ with

$$\begin{aligned}
\mathcal{K}^{(1)} &:= (U_\star^{(1)\top} \Delta_{U^{(1)}}, \mathbf{I}_{r_2}, \mathbf{I}_{r_3}) \cdot \mathcal{G}_\star, \\
\mathcal{K}^{(2)} &:= (\mathbf{I}_{r_1}, U_\star^{(2)\top} \Delta_{U^{(2)}}, \mathbf{I}_{r_3}) \cdot \mathcal{G}_\star, \\
\mathcal{K}^{(3)} &:= (\mathbf{I}_{r_1}, \mathbf{I}_{r_2}, U_\star^{(3)\top} \Delta_{U^{(3)}}) \cdot \mathcal{G}_\star.
\end{aligned}$$

As long as the choice of ϵ_0 is small enough such that $C_1 \epsilon_0 \rho^t \leq C_1 \epsilon_0 < 0.01$, the induction hypothesis (19a) tells us that.

$$\mathfrak{A}_{2,k} \geq \left\langle \mathcal{K}^{(k)}, \mathcal{K}^{(1)} + \mathcal{K}^{(2)} + \mathcal{K}^{(3)} \right\rangle - 0.01 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star). \tag{32}$$

- **Bounding \mathfrak{A}_3 .** This term is identical to the term that is bounded in [TMPB⁺22, Section B.2], which obeys

$$\mathfrak{A}_{3,k} \leq \left\| \mathcal{K}^{(1)} + \mathcal{K}^{(2)} + \mathcal{K}^{(3)} \right\|_{\mathbb{F}}^2 + C_2 \epsilon_0 \rho^t \text{dist}^2(\mathbf{F}_t, \mathbf{F}_\star)$$

for some constant $C_2 > 1$. As long as the choice of ϵ_0 is small enough such that $C_2 \epsilon_0 \rho^t \leq C_2 \epsilon_0 < 0.01$, the induction hypothesis (19a) results in

$$\mathfrak{A}_{3,k} \leq \left\| \mathcal{K}^{(1)} + \mathcal{K}^{(2)} + \mathcal{K}^{(3)} \right\|_{\mathbb{F}}^2 + 0.01 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star). \quad (33)$$

- **Bounding $\mathfrak{A}_{4,k}$.** To control $|\mathfrak{A}_{4,k}|$, we apply the definition of the matrix inner product to rewrite it as

$$\begin{aligned} |\mathfrak{A}_{4,k}| &= \left| \text{tr} \left(\mathcal{M}_k(\Delta_{\mathcal{S}}) \check{\mathbf{U}}^{(k)} (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} (\Sigma_\star^{(k)})^2 \Delta_{\mathbf{U}^{(k)}}^\top \right) \right| \\ &\leq \|\mathcal{M}_k(\Delta_{\mathcal{S}})\| \left\| \check{\mathbf{U}}^{(k)} (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} (\Sigma_\star^{(k)})^2 \Delta_{\mathbf{U}^{(k)}}^\top \right\|_* \\ &\leq 2\alpha \sqrt{n_1 n_2 n_3} \zeta_{t+1} \left\| \check{\mathbf{U}}^{(k)} (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} (\Sigma_\star^{(k)})^2 \Delta_{\mathbf{U}^{(k)}}^\top \right\|_*, \end{aligned}$$

where the second line follows from Hölder's inequality, and the third line follows from the bound of $\|\mathcal{M}_k(\Delta_{\mathcal{S}})\|$ from (30). To bound the remaining term, we have

$$\begin{aligned} \left\| \check{\mathbf{U}}^{(k)} (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} (\Sigma_\star^{(k)})^2 \Delta_{\mathbf{U}^{(k)}}^\top \right\|_* &\leq \sqrt{r_k} \left\| \check{\mathbf{U}}^{(k)} (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} (\Sigma_\star^{(k)})^2 \Delta_{\mathbf{U}^{(k)}}^\top \right\|_{\mathbb{F}} \\ &\leq \sqrt{r_k} \left\| \check{\mathbf{U}}^{(k)} (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} \Sigma_\star^{(k)} \right\| \left\| \Delta_{\mathbf{U}^{(k)}} \Sigma_\star^{(k)} \right\|_{\mathbb{F}} \\ &\leq \frac{\sqrt{r_k}}{(1 - \epsilon_0)^3} \left\| \Delta_{\mathbf{U}^{(k)}} \Sigma_\star^{(k)} \right\|_{\mathbb{F}}, \end{aligned}$$

where the last line follows from Lemma 8 (cf. (70c)). Plugging this into the previous inequality, we reach

$$|\mathfrak{A}_{4,k}| \leq \frac{2\alpha \zeta_{t+1} \sqrt{n_1 n_2 n_3 r_k}}{(1 - \epsilon_0)^3} \left\| \Delta_{\mathbf{U}^{(k)}} \Sigma_\star^{(k)} \right\|_{\mathbb{F}} \leq \frac{16\alpha \sqrt{\mu^3 r_1 r_2 r_3 r_k}}{(1 - \epsilon_0)^3} \rho^t \sigma_{\min}(\mathcal{X}_\star) \left\| \Delta_{\mathbf{U}^{(k)}} \Sigma_\star^{(k)} \right\|_{\mathbb{F}},$$

where the second inequality follows from the choice of ζ_{t+1} (cf. (22)). Finally, with the assumption on the sparsity level $\alpha \leq \frac{c_0 \epsilon_0}{\sqrt{\mu^3 r_1 r_2 r_3 r}}$ for sufficiently small c_0 and $\epsilon_0 < 0.01$, we have

$$|\mathfrak{A}_{4,k}| \leq \frac{16c_0}{(1 - \epsilon_0)^3} \left\| \Delta_{\mathbf{U}^{(k)}} \Sigma_\star^{(k)} \right\|_{\mathbb{F}} \epsilon_0 \rho^t \sigma_{\min}(\mathcal{X}_\star) \leq 0.15 \left\| \Delta_{\mathbf{U}^{(k)}} \Sigma_\star^{(k)} \right\|_{\mathbb{F}} \epsilon_0 \rho^t \sigma_{\min}(\mathcal{X}_\star). \quad (34)$$

- **Bounding $\mathfrak{A}_{5,k}$.** Similar to $\mathfrak{A}_{4,k}$, we first apply the definition of the matrix inner product and then Hölder's inequality, leading to

$$\begin{aligned} |\mathfrak{A}_{5,k}| &= \left| \text{tr} \left(\mathcal{M}_k(\Delta_{\mathcal{S}}) \check{\mathbf{U}}^{(k)} (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} (\Sigma_\star^{(k)})^2 (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} \check{\mathbf{U}}^{(k)\top} \Delta_{\check{\mathbf{U}}^{(k)}} \mathbf{U}_\star^{(k)\top} \right) \right| \\ &\leq \|\mathcal{M}_k(\Delta_{\mathcal{S}})\| \left\| \check{\mathbf{U}}^{(k)} (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} (\Sigma_\star^{(k)})^2 (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} \check{\mathbf{U}}^{(k)\top} \Delta_{\check{\mathbf{U}}^{(k)}} \mathbf{U}_\star^{(k)\top} \right\|_* \\ &\leq 2\zeta_{t+1} \alpha \sqrt{n_1 n_2 n_3 r_k} \left\| \check{\mathbf{U}}^{(k)} (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} (\Sigma_\star^{(k)})^2 (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} \check{\mathbf{U}}^{(k)\top} \Delta_{\check{\mathbf{U}}^{(k)}} \mathbf{U}_\star^{(k)\top} \right\|_{\mathbb{F}}, \end{aligned}$$

where the last line follows from (30), as well as the norm relation $\|\mathbf{A}\|_* \leq \sqrt{r_k} \|\mathbf{A}\|_{\mathbb{F}}$ for a matrix of rank at most r_k . To continue, noting that $\mathbf{U}_\star^{(k)}$ has orthonormal columns, we have

$$\begin{aligned} &\left\| \check{\mathbf{U}}^{(k)} (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} (\Sigma_\star^{(k)})^2 (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} \check{\mathbf{U}}^{(k)\top} \Delta_{\check{\mathbf{U}}^{(k)}} \mathbf{U}_\star^{(k)\top} \right\|_{\mathbb{F}} \\ &\leq \left\| \check{\mathbf{U}}^{(k)} (\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} \Sigma_\star^{(k)} \right\|_{\mathbb{F}}^2 \left\| \Delta_{\check{\mathbf{U}}^{(k)}} \right\|_{\mathbb{F}} \end{aligned}$$

$$\leq \frac{(1 + \epsilon_0 + \epsilon_0^2/3)}{(1 - \epsilon_0)^6} \left(\left\| \Delta_{\mathbf{U}^{(2)}} \Sigma_\star^{(2)} \right\|_{\mathbb{F}} + \left\| \Delta_{\mathbf{U}^{(3)}} \Sigma_\star^{(3)} \right\|_{\mathbb{F}} + \|\Delta_{\mathbf{g}}\|_{\mathbb{F}} \right) \leq \frac{2(1 + \epsilon_0 + \epsilon_0^2/3)}{(1 - \epsilon_0)^6} \text{dist}(\mathbf{F}_t, \mathbf{F}_\star),$$

where the penultimate line follows from Lemma 8 (cf. (70c) and (70e)), and the last line follows from Cauchy-Schwarz inequality. Plug this into the previous bound to arrive at

$$|\mathfrak{A}_{5,k}| \leq \frac{4\zeta_{t+1}\alpha\sqrt{n_1 n_2 n_3 r_k} (1 + \epsilon_0 + \epsilon_0^2/3)}{(1 - \epsilon_0)^6} \text{dist}(\mathbf{F}_t, \mathbf{F}_\star) \leq \frac{32\alpha\sqrt{\mu^3 r_1 r_2 r_3 r_k} (1 + \epsilon_0 + \epsilon_0^2/3)}{(1 - \epsilon_0)^6} \epsilon_0 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star),$$

where we have used $\text{dist}(\mathbf{F}_t, \mathbf{F}_\star) \leq \epsilon_0 \rho^t \sigma_{\min}(\mathcal{X}_\star)$ and the choice of ζ_{t+1} (cf. (22)). Finally, with the assumption on the sparsity level $\alpha \leq \frac{c_0 \epsilon_0}{\sqrt{\mu^3 r_1 r_2 r_3 r}}$ for sufficiently small c_0 and $\epsilon_0 < 0.01$, we have

$$|\mathfrak{A}_{5,k}| \leq \frac{32c_0 (1 + \epsilon_0 + \epsilon_0^2/3)}{(1 - \epsilon_0)^6} \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star) \leq 0.3 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star). \quad (35)$$

Putting things together. Summing (29) over all k , we obtain

$$\begin{aligned} & \sum_{k=1}^3 \left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)}) \Sigma_\star^{(k)} \right\|_{\mathbb{F}}^2 \\ &= (1 - \eta)^2 \sum_{k=1}^3 \left\| \Delta_{\mathbf{U}^{(k)}} \Sigma_\star^{(k)} \right\|_{\mathbb{F}}^2 + \eta^2 \sum_{k=1}^3 (\mathfrak{A}_{1,k} + \mathfrak{A}_{3,k} + 2\mathfrak{A}_{5,k}) - 2\eta(1 - \eta) \sum_{k=1}^3 (\mathfrak{A}_{2,k} + \mathfrak{A}_{4,k}). \end{aligned}$$

Plugging in our bounds in (31)-(35), we have

$$\begin{aligned} & \sum_{k=1}^3 \left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)}) \Sigma_\star^{(k)} \right\|_{\mathbb{F}}^2 \leq (1 - \eta)^2 \sum_{k=1}^3 \left\| \Delta_{\mathbf{U}^{(k)}} \Sigma_\star^{(k)} \right\|_{\mathbb{F}}^2 \\ & \quad + 2\eta^2 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star) - 2\eta(1 - \eta) \left\| \mathcal{K}^{(1)} + \mathcal{K}^{(2)} + \mathcal{K}^{(3)} \right\|_{\mathbb{F}}^2 + 0.06\eta(1 - \eta) \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star) \\ & \quad + 3\eta^2 \left\| \mathcal{K}^{(1)} + \mathcal{K}^{(2)} + \mathcal{K}^{(3)} \right\|_{\mathbb{F}}^2 + 0.15\eta(1 - \eta) \sum_{k=1}^3 \left\| \Delta_{\mathbf{U}^{(k)}} \Sigma_\star^{(k)} \right\|_{\mathbb{F}} \epsilon_0 \rho^t \sigma_{\min}(\mathcal{X}_\star). \end{aligned}$$

Note that when $0 < \eta < 2/5$,

$$-2\eta(1 - \eta) \left\| \mathcal{K}^{(1)} + \mathcal{K}^{(2)} + \mathcal{K}^{(3)} \right\|_{\mathbb{F}}^2 + 3\eta^2 \left\| \mathcal{K}^{(1)} + \mathcal{K}^{(2)} + \mathcal{K}^{(3)} \right\|_{\mathbb{F}}^2 = -\eta(2 - 5\eta) \left\| \mathcal{K}^{(1)} + \mathcal{K}^{(2)} + \mathcal{K}^{(3)} \right\|_{\mathbb{F}}^2 < 0.$$

Therefore, the previous bound can be simplified to

$$\begin{aligned} & \sum_{k=1}^3 \left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)}) \Sigma_\star^{(k)} \right\|_{\mathbb{F}}^2 \leq (1 - \eta)^2 \sum_{k=1}^3 \left\| \Delta_{\mathbf{U}^{(k)}} \Sigma_\star^{(k)} \right\|_{\mathbb{F}}^2 + 0.15\eta(1 - \eta) \sum_{k=1}^3 \left\| \Delta_{\mathbf{U}^{(k)}} \Sigma_\star^{(k)} \right\|_{\mathbb{F}} \epsilon_0 \rho^t \sigma_{\min}(\mathcal{X}_\star) \\ & \quad + 2\eta^2 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star) + 0.06\eta(1 - \eta) \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star). \end{aligned}$$

A.2 Proof of (26)

Taking the squared Frobenius norm of (25), it follows

$$\left\| ((\mathbf{Q}_t^{(1)})^{-1}, (\mathbf{Q}_t^{(2)})^{-1}, (\mathbf{Q}_t^{(3)})^{-1}) \cdot \mathbf{g}_{t+1} - \mathbf{g}_\star \right\|_{\mathbb{F}}^2 = (1 - \eta)^2 \|\Delta_{\mathbf{g}}\|_{\mathbb{F}}^2 - 2\eta(1 - \eta) \mathfrak{B}_1 + \eta^2 \mathfrak{B}_2, \quad (36)$$

where

$$\mathfrak{B}_1 = \left\langle \Delta_{\mathbf{g}}, \left((\mathbf{U}^{(1)\top} \mathbf{U}^{(1)})^{-1} \mathbf{U}^{(1)\top}, \dots, (\mathbf{U}^{(3)\top} \mathbf{U}^{(3)})^{-1} \mathbf{U}^{(3)\top} \right) \cdot \left((\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathbf{g}_\star - \mathcal{X}_\star + \Delta_{\mathbf{S}} \right) \right\rangle,$$

$$\mathfrak{B}_2 = \left\| \left((\mathbf{U}^{(1)\top} \mathbf{U}^{(1)})^{-1} \mathbf{U}^{(1)\top}, \dots, (\mathbf{U}^{(3)\top} \mathbf{U}^{(3)})^{-1} \mathbf{U}^{(3)\top} \right) \cdot \left((\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathbf{g}_\star - \mathcal{X}_\star + \Delta_{\mathbf{S}} \right) \right\|_{\mathbb{F}}^2.$$

We will now bound \mathfrak{B}_1 and \mathfrak{B}_2 separately.

Bounding \mathfrak{B}_1 . We start by breaking up the inner product into

$$\begin{aligned} \mathfrak{B}_1 &= \underbrace{\left\langle \Delta_{\mathcal{G}}, \left((\mathbf{U}^{(1)\top} \mathbf{U}^{(1)})^{-1} \mathbf{U}^{(1)\top}, \dots, (\mathbf{U}^{(3)\top} \mathbf{U}^{(3)})^{-1} \mathbf{U}^{(3)\top} \right) \cdot \left((\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathcal{G}_* - \mathcal{X}_* \right) \right\rangle}_{=:\mathfrak{B}_{1,1}} \\ &+ \underbrace{\left\langle \Delta_{\mathcal{G}}, \left((\mathbf{U}^{(1)\top} \mathbf{U}^{(1)})^{-1} \mathbf{U}^{(1)\top}, \dots, (\mathbf{U}^{(3)\top} \mathbf{U}^{(3)})^{-1} \mathbf{U}^{(3)\top} \right) \cdot \Delta_{\mathcal{S}} \right\rangle}_{=:\mathfrak{B}_{1,2}}. \end{aligned}$$

Note that $\mathfrak{B}_{1,1}$ is identical to the term that is bounded in [TMPB⁺22, Section B.3], which obeys

$$\mathfrak{B}_{1,1} \geq \sum_{k=1}^3 \left\| \mathbf{D}^{(k)} \right\|_{\mathbb{F}}^2 - C_1 \epsilon_0 \rho^t \text{dist}^2(\mathbf{F}_t, \mathbf{F}_*) \geq \sum_{k=1}^3 \left\| \mathbf{D}^{(k)} \right\|_{\mathbb{F}}^2 - 0.01 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_*), \quad (37)$$

where we have used the induction hypothesis (19a) and $C_1 \epsilon_0 \rho^t \leq C_1 \epsilon_0 < 0.01$ as long as ϵ_0 is small enough. Here,

$$\mathbf{D}^{(k)} = (\mathbf{U}^{(k)\top} \mathbf{U}^{(k)})^{-1/2} \mathbf{U}^{(k)\top} \Delta_{\mathbf{U}^{(k)} \Sigma_*^{(k)}}, \quad k = 1, 2, 3.$$

Turning to $\mathfrak{B}_{1,2}$, since the inner product is invariant to matricization, we flatten the tensor along the first mode to bound it as

$$\begin{aligned} |\mathfrak{B}_{1,2}| &= \left| \left\langle \mathcal{M}_1(\Delta_{\mathcal{G}}), (\mathbf{U}^{(1)\top} \mathbf{U}^{(1)})^{-1} \mathbf{U}^{(1)\top} \mathcal{M}_1(\Delta_{\mathcal{S}}) \left((\mathbf{U}^{(3)\top} \mathbf{U}^{(3)})^{-1} (\mathbf{U}^{(3)})^{\top} \otimes (\mathbf{U}^{(2)\top} \mathbf{U}^{(2)})^{-1} \mathbf{U}^{(2)\top} \right)^{\top} \right\rangle \right| \\ &\leq \|\mathcal{M}_1(\Delta_{\mathcal{G}})\|_* \left\| (\mathbf{U}^{(1)\top} \mathbf{U}^{(1)})^{-1} \mathbf{U}^{(1)\top} \mathcal{M}_1(\Delta_{\mathcal{S}}) \left((\mathbf{U}^{(3)\top} \mathbf{U}^{(3)})^{-1} \mathbf{U}^{(3)\top} \otimes (\mathbf{U}^{(2)\top} \mathbf{U}^{(2)})^{-1} \mathbf{U}^{(2)\top} \right)^{\top} \right\| \\ &\leq \sqrt{r_1} \|\Delta_{\mathcal{G}}\|_{\mathbb{F}} \prod_{k=1}^3 \left\| \mathbf{U}^{(k)} (\mathbf{U}^{(k)\top} \mathbf{U}^{(k)})^{-1} \right\| \|\mathcal{M}_1(\Delta_{\mathcal{S}})\|, \end{aligned}$$

where the second line uses Hölder's inequality, and the last line uses $\|\mathcal{M}_1(\Delta_{\mathcal{G}})\|_* \leq \sqrt{r_1} \|\Delta_{\mathcal{G}}\|_{\mathbb{F}}$ with the fact that $\mathcal{M}_1(\Delta_{\mathcal{G}})$ is of rank r_1 . To continue, invoke Lemma 8 (cf. (70b)) as well as (30) to further obtain

$$|\mathfrak{B}_{1,2}| \leq \frac{2\zeta_{t+1} \alpha \sqrt{n_1 n_2 n_3 r_1}}{(1 - \epsilon_0)^3} \|\Delta_{\mathcal{G}}\|_{\mathbb{F}} = \frac{16\alpha \sqrt{\mu^3 r_1^2 r_2 r_3}}{(1 - \epsilon_0)^3} \rho^t \sigma_{\min}(\mathcal{X}_*) \|\Delta_{\mathcal{G}}\|_{\mathbb{F}},$$

where the second equality follows from the choice of ζ_{t+1} (cf. (22)). Finally, with the assumption on the sparsity level $\alpha \leq \frac{c_0 \epsilon_0}{\sqrt{\mu^3 r_1 r_2 r_3 r}}$ for sufficiently small c_0 and $\epsilon_0 < 0.01$, we have

$$|\mathfrak{B}_{1,2}| \leq \frac{16c_0}{(1 - \epsilon_0)^3} \|\Delta_{\mathcal{G}}\|_{\mathbb{F}} \epsilon_0 \rho^t \sigma_{\min}(\mathcal{X}_*) \leq 0.15 \|\Delta_{\mathcal{G}}\|_{\mathbb{F}} \epsilon_0 \rho^t \sigma_{\min}(\mathcal{X}_*). \quad (38)$$

Put (37) and (38) together to see

$$\mathfrak{B}_1 \geq \sum_{k=1}^3 \left\| \mathbf{D}^{(k)} \right\|_{\mathbb{F}}^2 - 0.01 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_*) - 0.15 \|\Delta_{\mathcal{G}}\|_{\mathbb{F}} \epsilon_0 \rho^t \sigma_{\min}(\mathcal{X}_*). \quad (39)$$

Bounding \mathfrak{B}_2 . Expanding out the square and applying Cauchy-Schwarz, we can upper bound \mathfrak{B}_2 by

$$\begin{aligned} \mathfrak{B}_2 &\leq 2 \underbrace{\left\| \left((\mathbf{U}^{(1)\top} \mathbf{U}^{(1)})^{-1} \mathbf{U}^{(1)\top}, \dots, (\mathbf{U}^{(3)\top} \mathbf{U}^{(3)})^{-1} \mathbf{U}^{(3)\top} \right) \cdot \left((\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathcal{G}_* - \mathcal{X}_* \right) \right\|_{\mathbb{F}}^2}_{=:\mathfrak{B}_{2,1}} \\ &+ 2 \underbrace{\left\| \left((\mathbf{U}^{(1)\top} \mathbf{U}^{(1)})^{-1} \mathbf{U}^{(1)\top}, \dots, (\mathbf{U}^{(3)\top} \mathbf{U}^{(3)})^{-1} \mathbf{U}^{(3)\top} \right) \cdot \Delta_{\mathcal{S}} \right\|_{\mathbb{F}}^2}_{=:\mathfrak{B}_{2,2}}. \end{aligned}$$

$\mathfrak{B}_{2,1}$ is identical to the term that is bounded in [TMPB⁺22, Section B.4], which obeys

$$\mathfrak{B}_{2,1} \leq 3 \sum_{k=1}^3 \left\| \mathbf{D}^{(k)} \right\|_{\mathbb{F}}^2 + C_2 \epsilon_0 \rho^t \text{dist}^2(\mathbf{F}_t, \mathbf{F}_\star) \leq 3 \sum_{k=1}^3 \left\| \mathbf{D}^{(k)} \right\|_{\mathbb{F}}^2 + 0.01 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star), \quad (40)$$

where we make use of the induction hypothesis (19a) and $C_2 \epsilon_0 \rho^t \leq C_2 \epsilon_0 < 0.01$ as long as ϵ_0 is small enough.

For $\mathfrak{B}_{2,2}$, the matricization of the term in the norm along the first mode is of rank at most r_1 , so

$$\begin{aligned} \mathfrak{B}_{2,2} &= \left\| \left(((\mathbf{U}^{(1)})^\top \mathbf{U}^{(1)})^{-1} (\mathbf{U}^{(1)})^\top, ((\mathbf{U}^{(2)})^\top \mathbf{U}^{(2)})^{-1} (\mathbf{U}^{(2)})^\top, ((\mathbf{U}^{(3)})^\top \mathbf{U}^{(3)})^{-1} (\mathbf{U}^{(3)})^\top \right) \cdot \Delta \mathbf{s} \right\|_{\mathbb{F}}^2 \\ &\leq r_1 \left\| \left((\mathbf{U}^{(1)})^\top \mathbf{U}^{(1)} \right)^{-1} (\mathbf{U}^{(1)})^\top \mathcal{M}_1(\Delta \mathbf{s}) \left(((\mathbf{U}^{(3)})^\top \mathbf{U}^{(3)})^{-1} (\mathbf{U}^{(3)})^\top \otimes ((\mathbf{U}^{(2)})^\top \mathbf{U}^{(2)})^{-1} (\mathbf{U}^{(2)})^\top \right)^\top \right\|^2 \\ &\leq r_1 \left\| \left((\mathbf{U}^{(1)})^\top \mathbf{U}^{(1)} \right)^{-1} (\mathbf{U}^{(1)})^\top \right\|^2 \left\| \left((\mathbf{U}^{(2)})^\top \mathbf{U}^{(2)} \right)^{-1} (\mathbf{U}^{(2)})^\top \right\|^2 \left\| \left((\mathbf{U}^{(3)})^\top \mathbf{U}^{(3)} \right)^{-1} (\mathbf{U}^{(3)})^\top \right\|^2 \|\mathcal{M}_1(\Delta \mathbf{s})\|^2 \\ &\leq \frac{r_1}{(1 - \epsilon_0)^6} \|\mathcal{M}_1(\Delta \mathbf{s})\|^2, \end{aligned}$$

where the last line follows from Lemma 8 (cf. (70b)). To continue, we use (30) to obtain

$$\mathfrak{B}_{2,2} \leq \frac{4\zeta_{t+1}^2 \alpha^2 n_1 n_2 n_3 r_1}{(1 - \epsilon_0)^6} = \frac{128 \alpha^2 \mu^3 r_1^2 r_2 r_3}{(1 - \epsilon_0)^6} \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star),$$

where the second relation follows by the choice of ζ_{t+1} (cf. (22)). Lastly, with the assumption on the sparsity level $\alpha \leq \frac{c_0 \epsilon_0}{\sqrt{\mu^3 r_1 r_2 r_3}}$ for sufficiently small c_0 and $\epsilon_0 < 0.01$, we have

$$\mathfrak{B}_{2,2} \leq \frac{128 c_0^2}{(1 - \epsilon_0)^6} \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star) \leq 0.02 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star). \quad (41)$$

Combining (40) and (41), we get

$$\mathfrak{B}_2 \leq 6 \sum_{k=1}^3 \left\| \mathbf{D}^{(k)} \right\|_{\mathbb{F}}^2 + 0.06 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star). \quad (42)$$

Sum up. Going back to (36), we can substitute in our bounds for \mathfrak{B}_1 (cf. (39)) and \mathfrak{B}_2 (cf. (42)) to get

$$\begin{aligned} &\left\| (\mathbf{Q}_t^{(1)})^{-1}, (\mathbf{Q}_t^{(2)})^{-1}, (\mathbf{Q}_t^{(3)})^{-1} \right\|_{\mathbb{F}} \cdot \mathcal{G}_{t+1} - \mathcal{G}_\star \Big\|_{\mathbb{F}}^2 \\ &\leq (1 - \eta)^2 \|\Delta \mathcal{G}\|_{\mathbb{F}}^2 - 2\eta(1 - \eta) \left(\sum_{k=1}^3 \left\| \mathbf{D}^{(k)} \right\|_{\mathbb{F}}^2 - 0.01 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star) - 0.15 \|\Delta \mathcal{G}\|_{\mathbb{F}} \epsilon_0 \rho^t \sigma_{\min}(\mathcal{X}_\star) \right) \\ &\quad + \eta^2 \left(6 \sum_{k=1}^3 \left\| \mathbf{D}^{(k)} \right\|_{\mathbb{F}}^2 + 0.06 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star) \right). \end{aligned}$$

Notice that $-2\eta(1 - \eta) \|\mathbf{D}^{(k)}\|_{\mathbb{F}}^2 + 6\eta^2 \|\mathbf{D}^{(k)}\|_{\mathbb{F}}^2 = -2\eta(1 - 4\eta) \|\mathbf{D}^{(k)}\|_{\mathbb{F}}^2 \leq 0$ whenever $0 < \eta \leq 1/4$, leading to the conclusion that

$$\begin{aligned} \left\| (\mathbf{Q}_t^{(1)})^{-1}, (\mathbf{Q}_t^{(2)})^{-1}, (\mathbf{Q}_t^{(3)})^{-1} \right\|_{\mathbb{F}} \cdot \mathcal{G}_{t+1} - \mathcal{G}_\star \Big\|_{\mathbb{F}}^2 &\leq (1 - \eta)^2 \|\Delta \mathcal{G}\|_{\mathbb{F}}^2 + 2 \cdot 0.15 \eta (1 - \eta) \|\Delta \mathcal{G}\|_{\mathbb{F}} \epsilon_0 \rho^t \sigma_{\min}(\mathcal{X}_\star) \\ &\quad + 0.02 \eta (1 - \eta) \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star) + 0.06 \eta^2 \epsilon_0^2 \rho^{2t} \sigma_{\min}^2(\mathcal{X}_\star). \end{aligned}$$

B Proof of Lemma 2

In view of Lemma 1 and [TMPB⁺22, Lemma 6], the optimal alignment matrices $\{\mathbf{Q}_t^{(k)}\}_{k=1}^3$ (resp. $\{\mathbf{Q}_{t+1}^{(k)}\}_{k=1}^3$) between \mathbf{F}_t (resp. \mathbf{F}_{t+1}) and \mathbf{F}_\star exist. Fix any $k = 1, 2, 3$. By the triangle inequality, we have

$$\left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_{t+1}^{(k)} - \mathbf{U}_\star^{(k)}) \Sigma_\star^{(k)} \right\|_{2, \infty} \leq \left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)}) \Sigma_\star^{(k)} \right\|_{2, \infty} + \left\| \mathbf{U}_{t+1}^{(k)} (\mathbf{Q}_{t+1}^{(k)} - \mathbf{Q}_t^{(k)}) \Sigma_\star^{(k)} \right\|_{2, \infty}. \quad (43)$$

Below we control the two terms in turn.

Step 1: controlling $\left\| (U_{t+1}^{(k)} Q_t^{(k)} - U_\star^{(k)}) \Sigma_\star^{(k)} \right\|_{2,\infty}$. By the update rule, we have

$$(U_{t+1}^{(k)} Q_t^{(k)} - U_\star^{(k)}) \Sigma_\star^{(k)} = (1 - \eta) \Delta_{U^{(k)}} \Sigma_\star^{(k)} - \eta \left(\mathcal{M}_k(\Delta \mathcal{S}) + U_\star^{(k)} \Delta_{\check{U}^{(k)}}^\top \right) \check{U}^{(k)} (\check{U}^{(k)\top} \check{U}^{(k)})^{-1} \Sigma_\star^{(k)}.$$

Take the $\ell_{2,\infty}$ -norm of both sides and apply the triangle inequality to see that

$$\begin{aligned} \left\| (U_{t+1}^{(k)} Q_t^{(k)} - U_\star^{(k)}) \Sigma_\star^{(k)} \right\|_{2,\infty} &\leq (1 - \eta) \underbrace{\left\| \Delta_{U^{(k)}} \Sigma_\star^{(k)} \right\|_{2,\infty}}_{=: \mathfrak{C}_{1,k}} + \eta \underbrace{\left\| \mathcal{M}_k(\Delta \mathcal{S}) \check{U}^{(k)} (\check{U}^{(k)\top} \check{U}^{(k)})^{-1} \Sigma_\star^{(k)} \right\|_{2,\infty}}_{=: \mathfrak{C}_{2,k}} \\ &\quad + \eta \underbrace{\left\| U_\star^{(k)} \Delta_{\check{U}^{(k)}}^\top \check{U}^{(k)} (\check{U}^{(k)\top} \check{U}^{(k)})^{-1} \Sigma_\star^{(k)} \right\|_{2,\infty}}_{=: \mathfrak{C}_{3,k}}. \end{aligned}$$

We then proceed to bound each term separately.

- $\mathfrak{C}_{1,k}$ is captured by the induction hypothesis (19b), which directly implies

$$\mathfrak{C}_{1,k} \leq \rho^t \sqrt{\frac{\mu r_k}{n_k}} \sigma_{\min}(\mathcal{X}_\star). \quad (44)$$

- We now move on to $\mathfrak{C}_{2,k}$, which can be bounded by

$$\mathfrak{C}_{2,k} \leq \|\mathcal{M}_k(\Delta \mathcal{S})\|_{2,\infty} \left\| \check{U}^{(k)} (\check{U}^{(k)\top} \check{U}^{(k)})^{-1} \Sigma_\star^{(k)} \right\| \leq \frac{1}{(1 - \epsilon_0)^3} \|\mathcal{M}_k(\Delta \mathcal{S})\|_{2,\infty},$$

where the second inequality follows from Lemma 8 (cf. (70c)). Recall that $\Delta \mathcal{S}$ is α -sparse following the choice of ζ_{t+1} , which gives us

$$\|\mathcal{M}_k(\Delta \mathcal{S})\|_{2,\infty} \leq \sqrt{\frac{\alpha n_1 n_2 n_3}{n_k}} \|\Delta \mathcal{S}\|_\infty \leq 2 \sqrt{\frac{\alpha n_1 n_2 n_3}{n_k}} \zeta_{t+1} = 16 \sqrt{\frac{\alpha \mu^3 r_1 r_2 r_3}{n_k}} \rho^t \sigma_{\min}(\mathcal{X}_\star)$$

due to Lemma 11, (30), and the choice of ζ_{t+1} (cf. (22)). Plug this into the previous bound to obtain

$$\mathfrak{C}_{2,k} \leq \frac{16}{(1 - \epsilon_0)^3} \sqrt{\frac{\alpha \mu^3 r_1 r_2 r_3}{n_k}} \rho^t \sigma_{\min}(\mathcal{X}_\star) \leq 0.15 \sqrt{\frac{\mu r_k}{n_k}} \rho^t \sigma_{\min}(\mathcal{X}_\star), \quad (45)$$

where the last inequality follows from the assumption on the sparsity level $\alpha \leq \frac{c_1}{\mu^2 r_1 r_2 r_3}$ with a sufficiently small constant c_1 .

- Finally, for $\mathfrak{C}_{3,k}$, we have the upper bound

$$\begin{aligned} \mathfrak{C}_{3,k} &\leq \left\| U_\star^{(k)} \right\|_{2,\infty} \|\Delta_{\check{U}^{(k)}}\|_{\mathbb{F}} \left\| \check{U}^{(k)} (\check{U}^{(k)\top} \check{U}^{(k)})^{-1} \Sigma_\star^{(k)} \right\| \\ &\leq \sqrt{\frac{3\mu r_k}{n_k}} \left(1 + \epsilon_0 + \frac{\epsilon_0^2}{3} \right) \frac{1}{(1 - \epsilon_0)^3} \text{dist}(\mathbf{F}_t, \mathbf{F}_\star), \end{aligned}$$

where the second inequality follows from the incoherence assumption $\left\| U_\star^{(k)} \right\|_{2,\infty} \leq \sqrt{\frac{\mu r_k}{n_k}}$, and Lemma 8 (cf. (70c) and (70e)). Since $\text{dist}(\mathbf{F}_t, \mathbf{F}_\star) \leq \epsilon_0 \rho^t \sigma_{\min}(\mathcal{X}_\star)$, we arrive at

$$\mathfrak{C}_{3,k} \leq \frac{1}{(1 - \epsilon_0)^3} \sqrt{\frac{3\mu r_k}{n_k}} \left(1 + \epsilon_0 + \frac{\epsilon_0^2}{3} \right) \epsilon_0 \rho^t \sigma_{\min}(\mathcal{X}_\star) \leq 0.02 \sqrt{\frac{\mu r_k}{n_k}} \rho^t \sigma_{\min}(\mathcal{X}_\star) \quad (46)$$

as long as $\epsilon_0 < 0.01$.

Combining (44), (45), and (46) together, we reach the conclusion that

$$\left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\|_{2,\infty} \leq (1 - 0.83\eta) \sqrt{\frac{\mu r_k}{n_k}} \rho^t \sigma_{\min}(\boldsymbol{\mathcal{X}}_\star). \quad (47)$$

In view of the basic relation

$$\left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\|_{2,\infty} \geq \left\| \mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)} \right\|_{2,\infty} \sigma_{\min}(\boldsymbol{\mathcal{X}}_\star),$$

this also implies

$$\left\| \mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)} \right\|_{2,\infty} \leq (1 - 0.83\eta) \sqrt{\frac{\mu r_k}{n_k}} \rho^t. \quad (48)$$

Step 2: controlling $\left\| \mathbf{U}_{t+1}^{(k)} (\mathbf{Q}_{t+1}^{(k)} - \mathbf{Q}_t^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\|_{2,\infty}$. Observe that

$$\begin{aligned} \left\| \mathbf{U}_{t+1}^{(k)} (\mathbf{Q}_{t+1}^{(k)} - \mathbf{Q}_t^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\|_{2,\infty} &= \left\| \mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} (\mathbf{Q}_t^{(k)})^{-1} (\mathbf{Q}_{t+1}^{(k)} - \mathbf{Q}_t^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\|_{2,\infty} \\ &\leq \left\| \mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} \right\|_{2,\infty} \left\| (\mathbf{Q}_t^{(k)})^{-1} (\mathbf{Q}_{t+1}^{(k)} - \mathbf{Q}_t^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\|. \end{aligned}$$

For the first term $\left\| \mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} \right\|_{2,\infty}$, we have

$$\begin{aligned} \left\| \mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} \right\|_{2,\infty} &\leq \left\| \mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)} \right\|_{2,\infty} + \left\| \mathbf{U}_\star^{(k)} \right\|_{2,\infty} \\ &\leq ((1 - 0.83\eta)\rho^t + 1) \sqrt{\frac{\mu r_k}{n_k}} \leq (2 - 0.83\eta) \sqrt{\frac{\mu r_k}{n_k}}, \end{aligned}$$

where the second line follows from (48) and the incoherence assumption $\left\| \mathbf{U}_\star^{(k)} \right\|_{2,\infty} \leq \sqrt{\frac{\mu r_k}{n_k}}$.

Moving on to the second term $\left\| (\mathbf{Q}_t^{(k)})^{-1} (\mathbf{Q}_{t+1}^{(k)} - \mathbf{Q}_t^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\|$, we plan to invoke Lemma 7. Given the assumption of the sparsity level α satisfies the requirement of Lemma 1, following Lemma 1 as well as its proof (cf. (20)), we have

$$\begin{aligned} \left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_{t+1}^{(k)} - \mathbf{U}_\star^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\|_{\text{F}} &\leq \text{dist}(\mathbf{F}_{t+1}, \mathbf{F}_\star) \leq \epsilon_0 \rho^{t+1} \sigma_{\min}(\boldsymbol{\mathcal{X}}_\star), \\ \left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\|_{\text{F}} &\leq \epsilon_0 \rho^{t+1} \sigma_{\min}(\boldsymbol{\mathcal{X}}_\star), \end{aligned}$$

which in turn implies

$$\max \left\{ \left\| \mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)} \right\|, \frac{\left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_{t+1}^{(k)} - \mathbf{U}_\star^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\|}{\sigma_{\min}(\boldsymbol{\mathcal{X}}_\star)}, \frac{\left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\|}{\sigma_{\min}(\boldsymbol{\mathcal{X}}_\star)} \right\} \leq \epsilon_0 \rho^{t+1} < 1. \quad (49)$$

Setting $\mathbf{U} := \mathbf{U}_{t+1}^{(k)}$, $\mathbf{U}_\star := \mathbf{U}_\star^{(k)}$, $\mathbf{Q} := \mathbf{Q}_{t+1}^{(k)}$, $\bar{\mathbf{Q}} := \mathbf{Q}_t^{(k)}$ and $\boldsymbol{\Sigma} := \boldsymbol{\Sigma}_\star^{(k)}$, (49) demonstrates that Lemma 7 is applicable, leading to

$$\begin{aligned} \left\| (\mathbf{Q}_t^{(k)})^{-1} (\mathbf{Q}_{t+1}^{(k)} - \mathbf{Q}_t^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\| &\leq \frac{\left\| \mathbf{U}_{t+1}^{(k)} (\mathbf{Q}_{t+1}^{(k)} - \mathbf{Q}_t^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\|}{\sigma_{\min}(\mathbf{U}_\star) - \left\| \mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)} \right\|} \\ &\leq \frac{\left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_{t+1}^{(k)} - \mathbf{U}_\star^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\| + \left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_t^{(k)} - \mathbf{U}_\star^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\|}{1 - \epsilon_0} \\ &\leq \frac{2\epsilon_0}{1 - \epsilon_0} \rho^{t+1} \sigma_{\min}(\boldsymbol{\mathcal{X}}_\star), \end{aligned}$$

where we used (49) in the second and third inequalities. Combining the above two bounds, we have

$$\left\| \mathbf{U}_{t+1}^{(k)} (\mathbf{Q}_{t+1}^{(k)} - \mathbf{Q}_t^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\|_{2,\infty} \leq (2 - 0.83\eta) \frac{2\epsilon_0}{1 - \epsilon_0} \sqrt{\frac{\mu r_k}{n_k}} \rho^{t+1} \sigma_{\min}(\boldsymbol{\mathcal{X}}_\star). \quad (50)$$

Step 3: combining the bounds. Plug (47) and (50) into (43) to get

$$\begin{aligned}
\left\| (\mathbf{U}_{t+1}^{(k)} \mathbf{Q}_{t+1}^{(k)} - \mathbf{U}_\star^{(k)}) \boldsymbol{\Sigma}_\star^{(k)} \right\|_{2,\infty} &\leq \left[(1 - 0.83\eta) + (2 - 0.83\eta) \frac{2\rho\epsilon_0}{1 - \epsilon_0} \right] \sqrt{\frac{\mu r_k}{n_k}} \rho^t \sigma_{\min}(\boldsymbol{\mathcal{X}}_\star) \\
&\leq (1.05 - 0.84\eta) \sqrt{\frac{\mu r_k}{n_k}} \rho^t \sigma_{\min}(\boldsymbol{\mathcal{X}}_\star) \\
&\leq (1 - 0.45\eta) \sqrt{\frac{\mu r_k}{n_k}} \rho^t \sigma_{\min}(\boldsymbol{\mathcal{X}}_\star) = \sqrt{\frac{\mu r_k}{n_k}} \rho^{t+1} \sigma_{\min}(\boldsymbol{\mathcal{X}}_\star)
\end{aligned}$$

where the last line follows from $\eta \geq \frac{1}{7}$ and $\rho = 1 - 0.45\eta$.

C Proof of Lemma 3

In view of [TMPB⁺22, Lemma 8], one has

$$\text{dist}(\mathbf{F}_0, \mathbf{F}_\star) \leq (\sqrt{2} + 1)^{\frac{3}{2}} \left\| (\mathbf{U}_0^{(1)}, \mathbf{U}_0^{(2)}, \mathbf{U}_0^{(3)}) \cdot \boldsymbol{\mathcal{G}}_0 - \boldsymbol{\mathcal{X}}_\star \right\|_{\text{F}}, \quad (51)$$

where we have used the definition $\mathbf{F}_0 = (\mathbf{U}_0^{(1)}, \mathbf{U}_0^{(2)}, \mathbf{U}_0^{(3)}) \cdot \boldsymbol{\mathcal{G}}_0$. As a result, we focus on bounding the term $\left\| (\mathbf{U}_0^{(1)}, \mathbf{U}_0^{(2)}, \mathbf{U}_0^{(3)}) \cdot \boldsymbol{\mathcal{G}}_0 - \boldsymbol{\mathcal{X}}_\star \right\|_{\text{F}}$ below.

Recall from the definition of the HOSVD that $\boldsymbol{\mathcal{G}}_0 = ((\mathbf{U}_0^{(1)})^\top, (\mathbf{U}_0^{(2)})^\top, (\mathbf{U}_0^{(3)})^\top) \cdot (\boldsymbol{\mathcal{Y}} - \mathcal{T}_{\zeta_0}(\boldsymbol{\mathcal{Y}}))$, which in turn implies

$$(\mathbf{U}_0^{(1)}, \mathbf{U}_0^{(2)}, \mathbf{U}_0^{(3)}) \cdot \boldsymbol{\mathcal{G}}_0 = \left(\mathbf{U}_0^{(1)} \mathbf{U}_0^{(1)\top}, \mathbf{U}_0^{(2)} \mathbf{U}_0^{(2)\top}, \mathbf{U}_0^{(3)} \mathbf{U}_0^{(3)\top} \right) \cdot (\boldsymbol{\mathcal{Y}} - \mathcal{T}_{\zeta_0}(\boldsymbol{\mathcal{Y}})). \quad (52)$$

Note that $\mathbf{U}_0^{(k)}$ has orthonormal columns. We thus define $\mathbf{P}^{(k)} := \mathbf{U}_0^{(k)} \mathbf{U}_0^{(k)\top}$, the orthogonal projection onto the column space of $\mathbf{U}_0^{(k)}$. This allows us to rewrite the squared Frobenius norm in (51) as

$$\left\| (\mathbf{U}_0^{(1)}, \mathbf{U}_0^{(2)}, \mathbf{U}_0^{(3)}) \cdot \boldsymbol{\mathcal{G}}_0 - \boldsymbol{\mathcal{X}}_\star \right\|_{\text{F}}^2 = \left\| (\mathbf{P}^{(1)}, \mathbf{P}^{(2)}, \mathbf{P}^{(3)}) \cdot (\boldsymbol{\mathcal{Y}} - \mathcal{T}_{\zeta_0}(\boldsymbol{\mathcal{Y}})) - \boldsymbol{\mathcal{X}}_\star \right\|_{\text{F}}^2.$$

Since $\boldsymbol{\mathcal{X}}_\star$ can be decomposed into a sum of orthogonal projections and its orthogonal complements, namely,

$$\begin{aligned}
\boldsymbol{\mathcal{X}}_\star &= (\mathbf{P}^{(1)}, \mathbf{P}^{(2)}, \mathbf{P}^{(3)}) \cdot \boldsymbol{\mathcal{X}}_\star + (\mathbf{I}_{n_1} - \mathbf{P}^{(1)}, \mathbf{P}^{(2)}, \mathbf{P}^{(3)}) \cdot \boldsymbol{\mathcal{X}}_\star \\
&\quad + (\mathbf{I}_{n_1}, \mathbf{I}_{n_2} - \mathbf{P}^{(2)}, \mathbf{P}^{(3)}) \cdot \boldsymbol{\mathcal{X}}_\star + (\mathbf{I}_{n_1}, \mathbf{I}_{n_2}, \mathbf{I}_{n_3} - \mathbf{P}^{(3)}) \cdot \boldsymbol{\mathcal{X}}_\star,
\end{aligned}$$

we have the following identity

$$\begin{aligned}
\left\| (\mathbf{U}_0^{(1)}, \mathbf{U}_0^{(2)}, \mathbf{U}_0^{(3)}) \cdot \boldsymbol{\mathcal{G}}_0 - \boldsymbol{\mathcal{X}}_\star \right\|_{\text{F}}^2 &= \left\| (\mathbf{P}^{(1)}, \mathbf{P}^{(2)}, \mathbf{P}^{(3)}) \cdot (\boldsymbol{\mathcal{Y}} - \mathcal{T}_{\zeta_0}(\boldsymbol{\mathcal{Y}})) - \boldsymbol{\mathcal{X}}_\star \right\|_{\text{F}}^2 \\
&+ \left\| (\mathbf{I}_{n_1} - \mathbf{P}^{(1)}, \mathbf{P}^{(2)}, \mathbf{P}^{(3)}) \cdot \boldsymbol{\mathcal{X}}_\star \right\|_{\text{F}}^2 + \left\| (\mathbf{I}_{n_1}, \mathbf{I}_{n_2} - \mathbf{P}^{(2)}, \mathbf{P}^{(3)}) \cdot \boldsymbol{\mathcal{X}}_\star \right\|_{\text{F}}^2 + \left\| (\mathbf{I}_{n_1}, \mathbf{I}_{n_2}, \mathbf{I}_{n_3} - \mathbf{P}^{(3)}) \cdot \boldsymbol{\mathcal{X}}_\star \right\|_{\text{F}}^2.
\end{aligned} \quad (53)$$

In what follows, we bound each term respectively.

Bounding the first term. Matricize along the first mode and change to the operator norm to obtain

$$\begin{aligned}
\left\| (\mathbf{P}^{(1)}, \mathbf{P}^{(2)}, \mathbf{P}^{(3)}) \cdot (\boldsymbol{\mathcal{Y}} - \mathcal{T}_{\zeta_0}(\boldsymbol{\mathcal{Y}})) - \boldsymbol{\mathcal{X}}_\star \right\|_{\text{F}} &\leq \sqrt{r_1} \left\| \mathbf{P}^{(1)} \mathcal{M}_1 (\boldsymbol{\mathcal{Y}} - \mathcal{T}_{\zeta_0}(\boldsymbol{\mathcal{Y}})) - \boldsymbol{\mathcal{X}}_\star \right\| \\
&\leq \sqrt{r_1} \left\| \mathcal{M}_1 (\boldsymbol{\mathcal{Y}} - \mathcal{T}_{\zeta_0}(\boldsymbol{\mathcal{Y}})) - \boldsymbol{\mathcal{X}}_\star \right\| \\
&= \sqrt{r_1} \left\| \mathcal{M}_1 (\boldsymbol{\mathcal{S}}_\star - \mathcal{T}_{\zeta_0}(\boldsymbol{\mathcal{Y}})) \right\|,
\end{aligned} \quad (54)$$

where the last relation holds due to the definition of $\boldsymbol{\mathcal{Y}}$.

Bounding the remaining three terms. We present the bound on the second term, while the remaining two can be bounded in a similar fashion. Matricize along the first mode, and in view of the fact that $\|\mathbf{P}^{(k)}\| \leq 1$, we have

$$\begin{aligned} \left\| (\mathbf{I}_{n_1} - \mathbf{P}^{(1)}, \mathbf{P}^{(2)}, \mathbf{P}^{(3)}) \cdot \mathcal{X}_* \right\|_{\mathbb{F}} &\leq \left\| (\mathbf{I}_{n_1} - \mathbf{P}^{(1)}) \mathcal{M}_1(\mathcal{X}_*) \right\|_{\mathbb{F}} \\ &\leq \sqrt{r_1} \left\| (\mathbf{I}_{n_1} - \mathbf{P}^{(1)}) \mathcal{M}_1(\mathcal{X}_*) \right\| \\ &= \sqrt{r_1} \left\| (\mathbf{I}_{n_1} - \mathbf{P}^{(1)}) \mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y}) - \mathcal{S}_* + \mathcal{T}_{\zeta_0}(\mathcal{Y})) \right\| \\ &\leq \sqrt{r_1} \|\mathcal{M}_1(\mathcal{S}_* - \mathcal{T}_{\zeta_0}(\mathcal{Y}))\| + \sqrt{r_1} \left\| (\mathbf{I}_{n_1} - \mathbf{P}^{(1)}) \mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y})) \right\|. \end{aligned}$$

To continue, note that $\mathbf{P}^{(1)} \mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y}))$ is the best rank- r_1 approximation to $\mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y}))$, which implies

$$\begin{aligned} \left\| (\mathbf{I}_{n_1} - \mathbf{P}^{(1)}) \mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y})) \right\| &\leq \sigma_{r_1+1}(\mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y}))) \\ &\leq \sigma_{r_1+1}(\mathcal{M}_1(\mathcal{X}_*)) + \|\mathcal{M}_1(\mathcal{S}_* - \mathcal{T}_{\zeta_0}(\mathcal{Y}))\| = \|\mathcal{M}_1(\mathcal{S}_* - \mathcal{T}_{\zeta_0}(\mathcal{Y}))\|, \end{aligned}$$

where the last line follows from Weyl's inequality and the fact that $\mathcal{M}_1(\mathcal{X}_*)$ has rank r_1 . Plug this into the previous inequality to obtain

$$\left\| (\mathbf{I}_{n_1} - \mathbf{P}^{(1)}, \mathbf{P}^{(2)}, \mathbf{P}^{(3)}) \cdot \mathcal{X}_* \right\|_{\mathbb{F}} \leq 2\sqrt{r_1} \|\mathcal{M}_1(\mathcal{S}_* - \mathcal{T}_{\zeta_0}(\mathcal{Y}))\|. \quad (55)$$

Plug our bounds (54) and (55) into (53) to obtain

$$\begin{aligned} \left\| (\mathbf{U}_0^{(1)}, \mathbf{U}_0^{(2)}, \mathbf{U}_0^{(3)}) \cdot \mathcal{G}_0 - \mathcal{X}_* \right\|_{\mathbb{F}}^2 &\leq r_1 \|\mathcal{M}_1(\mathcal{S}_* - \mathcal{T}_{\zeta_0}(\mathcal{Y}))\|^2 + 4r_1 \|\mathcal{M}_1(\mathcal{S}_* - \mathcal{T}_{\zeta_0}(\mathcal{Y}))\|^2 \\ &\quad + 4r_2 \|\mathcal{M}_2(\mathcal{S}_* - \mathcal{T}_{\zeta_0}(\mathcal{Y}))\|^2 + 4r_3 \|\mathcal{M}_3(\mathcal{S}_* - \mathcal{T}_{\zeta_0}(\mathcal{Y}))\|^2. \end{aligned} \quad (56)$$

It then boils down to controlling $\|\mathcal{M}_k(\mathcal{S}_* - \mathcal{T}_{\zeta_0}(\mathcal{Y}))\|$ for $k = 1, 2, 3$. With our choice of ζ_0 (i.e. $\|\mathcal{X}_*\|_{\infty} \leq \zeta_0 \leq 2\|\mathcal{X}_*\|_{\infty}$), setting $\mathcal{X} = \mathbf{0}$ in Lemma 12 guarantees that $\mathcal{M}_k(\mathcal{S}_* - \mathcal{T}_{\zeta_0}(\mathcal{Y}))$ is α -sparse for all k . Hence, we can apply Lemma 11 to arrive at

$$\|\mathcal{M}_k(\mathcal{S}_* - \mathcal{T}_{\zeta_0}(\mathcal{Y}))\| \leq \alpha \sqrt{n_1 n_2 n_3} \|\mathcal{S}_* - \mathcal{T}_{\zeta_0}(\mathcal{Y})\|_{\infty} \leq 2\alpha \sqrt{n_1 n_2 n_3} \zeta_0 \leq 4\alpha \sqrt{n_1 n_2 n_3} \|\mathcal{X}_*\|_{\infty}, \quad (57)$$

where the penultimate inequality follows from Lemma 12 (cf. (76)). Plug (57) into (56) to obtain

$$\left\| (\mathbf{U}_0^{(1)}, \mathbf{U}_0^{(2)}, \mathbf{U}_0^{(3)}) \cdot \mathcal{G}_0 - \mathcal{X}_* \right\|_{\mathbb{F}}^2 \leq 208\alpha^2 n_1 n_2 n_3 r \|\mathcal{X}_*\|_{\infty}^2 \leq 208\alpha^2 \mu^3 r_1 r_2 r_3 r \kappa^2 \sigma_{\min}^2(\mathcal{X}_*),$$

where the second inequality follows from Lemma 5. Under the assumption that $\alpha \leq \frac{c_0}{\sqrt{\mu^3 r_1 r_2 r_3 r \kappa}}$, it follows

$$\left\| (\mathbf{U}_0^{(1)}, \mathbf{U}_0^{(2)}, \mathbf{U}_0^{(3)}) \cdot \mathcal{G}_0 - \mathcal{X}_* \right\|_{\mathbb{F}}^2 \leq 208c_0^2 \sigma_{\min}^2(\mathcal{X}_*).$$

This combined with (51) finishes the proof.

D Proof of Lemma 4

We provide the control on the first mode, as the other two modes can be bounded using the same arguments.

We begin with a useful decomposition of the quantity we care about, whose proof will be supplied in the end of this section:

$$\begin{aligned} \Delta_{\mathbf{U}^{(1)}} \Sigma_*^{(1)} &= \mathcal{M}_1(\mathcal{S}_* - \mathcal{T}_{\zeta_0}(\mathcal{Y})) \mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y}))^{\top} \mathbf{U}_0^{(1)} (\mathbf{Q}_0^{(1)})^{-\top} (\check{\mathbf{U}}^{(1)\top} \check{\mathbf{U}}^{(1)})^{-1} \Sigma_*^{(1)} \\ &\quad + \mathbf{U}_*^{(1)} \Delta_{\check{\mathbf{U}}^{(1)}}^{\top} \mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y}))^{\top} \mathbf{U}_0^{(1)} (\mathbf{Q}_0^{(1)})^{-\top} (\check{\mathbf{U}}^{(1)\top} \check{\mathbf{U}}^{(1)})^{-1} \Sigma_*^{(1)}. \end{aligned} \quad (58)$$

Taking the $\ell_{2,\infty}$ -norm and using the triangle inequality, we obtain

$$\begin{aligned} \left\| \mathbf{\Delta}_{\check{U}^{(1)}} \mathbf{\Sigma}_*^{(1)} \right\|_{2,\infty} &\leq \underbrace{\left\| \mathbf{U}_*^{(1)} \mathbf{\Delta}_{\check{U}^{(1)}}^\top \mathcal{M}_1 (\mathbf{Y} - \mathcal{T}_{\zeta_0}(\mathbf{Y}))^\top \mathbf{U}_0^{(1)} (\mathbf{Q}_0^{(1)})^{-\top} (\check{U}^{(1)\top} \check{U}^{(1)})^{-1} \mathbf{\Sigma}_*^{(1)} \right\|_{2,\infty}}_{=: \mathfrak{A}_1} \\ &\quad + \underbrace{\left\| \mathcal{M}_1 (\mathbf{s}_* - \mathcal{T}_{\zeta_0}(\mathbf{Y})) \mathcal{M}_1 (\mathbf{Y} - \mathcal{T}_{\zeta_0}(\mathbf{Y}))^\top \mathbf{U}_0^{(1)} (\mathbf{Q}_0^{(1)})^{-\top} (\check{U}^{(1)\top} \check{U}^{(1)})^{-1} \mathbf{\Sigma}_*^{(1)} \right\|_{2,\infty}}_{=: \mathfrak{A}_2}. \end{aligned} \quad (59)$$

We now proceed to bound these two terms separately.

Step 1: bounding \mathfrak{A}_1 . To begin, note that

$$\begin{aligned} \mathfrak{A}_1 &\leq \left\| \mathbf{U}_*^{(1)} \right\|_{2,\infty} \left\| \mathbf{\Delta}_{\check{U}^{(1)}} \right\| \left\| \mathcal{M}_1 (\mathbf{Y} - \mathcal{T}_{\zeta_0}(\mathbf{Y}))^\top \mathbf{U}_0^{(1)} (\mathbf{Q}_0^{(1)})^{-\top} (\check{U}^{(1)\top} \check{U}^{(1)})^{-1} \mathbf{\Sigma}_*^{(1)} \right\| \\ &= \left\| \mathbf{U}_*^{(1)} \right\|_{2,\infty} \left\| \mathbf{\Delta}_{\check{U}^{(1)}} \right\| \left\| \check{U}^{(1)} (\check{U}^{(1)\top} \check{U}^{(1)})^{-1} \mathbf{\Sigma}_*^{(1)} \right\|, \end{aligned}$$

where in the second line we have used the relation

$$\check{U}_0^{(1)\top} \check{U}_0^{(1)} = \mathbf{U}_0^{(1)\top} \mathcal{M}_1 (\mathbf{Y} - \mathcal{T}_{\zeta_0}(\mathbf{Y})) \mathcal{M}_1 (\mathbf{Y} - \mathcal{T}_{\zeta_0}(\mathbf{Y}))^\top \mathbf{U}_0^{(1)}$$

given by Lemma 6 (cf. (66)), and the short-hand notation in (21). Invoking Lemma 8 (cf. (70c)) and the incoherence assumption $\left\| \mathbf{U}_*^{(1)} \right\|_{2,\infty} \leq \sqrt{\frac{\mu r_1}{n_1}}$, we arrive at

$$\mathfrak{A}_1 \leq \frac{1}{(1 - \epsilon_0)^3} \sqrt{\frac{\mu r_1}{n_1}} \left\| \mathbf{\Delta}_{\check{U}^{(1)}} \right\|.$$

Furthermore, by Lemma 8 (cf. (70e)), it holds that

$$\left\| \mathbf{\Delta}_{\check{U}^{(1)}} \right\| \leq \left\| \mathbf{\Delta}_{\check{U}^{(1)}} \right\|_{\mathbf{F}} \leq 2 \left(1 + \epsilon_0 + \frac{\epsilon_0^2}{3} \right) \text{dist}(\mathbf{F}_0, \mathbf{F}_*) \leq 2 \left(1 + \epsilon_0 + \frac{\epsilon_0^2}{3} \right) \epsilon_0 \sigma_{\min}(\mathbf{X}_*),$$

leading to

$$\mathfrak{A}_1 \leq \frac{2\epsilon_0}{(1 - \epsilon_0)^3} \left(1 + \epsilon_0 + \frac{\epsilon_0^2}{3} \right) \sqrt{\frac{\mu r_1}{n_1}} \sigma_{\min}(\mathbf{X}_*) \leq 0.57 \sqrt{\frac{\mu r_1}{n_1}} \sigma_{\min}(\mathbf{X}_*), \quad (60)$$

where the last inequality holds as long as $\epsilon_0 \leq 0.15$.

Step 2: bounding \mathfrak{A}_2 . For \mathfrak{A}_2 , we have

$$\begin{aligned} \mathfrak{A}_2 &\leq \left\| \mathcal{M}_1 (\mathbf{s}_* - \mathcal{T}_{\zeta_0}(\mathbf{Y})) \right\|_{1,\infty} \left\| \mathcal{M}_1 (\mathbf{Y} - \mathcal{T}_{\zeta_0}(\mathbf{Y}))^\top \mathbf{U}_0^{(1)} (\mathbf{Q}_0^{(1)})^{-\top} (\mathbf{\Sigma}_*^{(1)})^{-1} \right\|_{2,\infty} \left\| \mathbf{\Sigma}_*^{(1)} (\check{U}^{(1)\top} \check{U}^{(1)})^{-1} \mathbf{\Sigma}_*^{(1)} \right\| \\ &\leq \frac{\alpha n_2 n_3}{(1 - \epsilon_0)^6} \left\| \mathbf{s}_* - \mathcal{T}_{\zeta_0}(\mathbf{X}_* + \mathbf{s}_*) \right\|_{\infty} \left\| \mathcal{M}_1 (\mathbf{Y} - \mathcal{T}_{\zeta_0}(\mathbf{Y}))^\top \mathbf{U}_0^{(1)} (\mathbf{Q}_0^{(1)})^{-\top} (\mathbf{\Sigma}_*^{(1)})^{-1} \right\|_{2,\infty}. \end{aligned}$$

where we have used the α -sparsity of $\mathcal{M}_1 (\mathbf{s}_* - \mathcal{T}_{\zeta_0}(\mathbf{Y}))$ given by Lemma 12 and Lemma 8 (cf. (70d)) in the second line. Apply Lemma 12 with $\mathbf{X} = \mathbf{0}$ to get

$$\begin{aligned} \mathfrak{A}_2 &\leq \frac{2\alpha n_2 n_3}{(1 - \epsilon_0)^6} \zeta_0 \left\| \mathcal{M}_1 (\mathbf{Y} - \mathcal{T}_{\zeta_0}(\mathbf{Y}))^\top \mathbf{U}_0^{(1)} (\mathbf{Q}_0^{(1)})^{-\top} (\mathbf{\Sigma}_*^{(1)})^{-1} \right\|_{2,\infty} \\ &\leq \frac{4c_0}{(1 - \epsilon_0)^6} \sqrt{\frac{n_2 n_3}{\mu n_1 r_1 r_2 r_3}} \sigma_{\min}(\mathbf{X}_*) \left\| \mathcal{M}_1 (\mathbf{Y} - \mathcal{T}_{\zeta_0}(\mathbf{Y}))^\top \mathbf{U}_0^{(1)} (\mathbf{Q}_0^{(1)})^{-\top} (\mathbf{\Sigma}_*^{(1)})^{-1} \right\|_{2,\infty}, \end{aligned} \quad (61)$$

where the second line follows from $\zeta_0 \leq 2 \|\mathbf{X}_*\|_{\infty} \leq 2 \sqrt{\frac{\mu^3 r_1 r_2 r_3}{n_1 n_2 n_3}} \kappa \sigma_{\min}(\mathbf{X}_*)$ (cf. Lemma 5), as well as the assumption $\alpha \leq \frac{c_0}{\mu^2 r_1 r_2 r_3 \kappa}$. To continue,

$$\left\| \mathcal{M}_1 (\mathbf{Y} - \mathcal{T}_{\zeta_0}(\mathbf{Y}))^\top \mathbf{U}_0^{(1)} (\mathbf{Q}_0^{(1)})^{-\top} (\mathbf{\Sigma}_*^{(1)})^{-1} \right\|_{2,\infty} \leq r_1 \left\| \mathcal{M}_1 (\mathbf{Y} - \mathcal{T}_{\zeta_0}(\mathbf{Y}))^\top \mathbf{U}_0^{(1)} (\mathbf{Q}_0^{(1)})^{-\top} (\mathbf{\Sigma}_*^{(1)})^{-1} \right\|_{\infty}$$

$$\begin{aligned}
&\leq r_1 \left\| (\boldsymbol{\Sigma}_\star^{(1)})^{-1} (\mathbf{Q}_0^{(1)})^{-1} \mathbf{U}_0^{(1)\top} \mathcal{M}_1 (\mathbf{y} - \mathcal{T}_{\zeta_0}(\mathbf{y})) \right\|_{2,\infty} \\
&= r_1 \left\| (\boldsymbol{\Sigma}_\star^{(1)})^{-1} \check{\mathbf{U}}^{(1)\top} \check{\mathbf{U}}^{(1)} (\boldsymbol{\Sigma}_\star^{(1)})^{-1} \right\|_{\infty}^{1/2} \\
&= r_1 \left\| \check{\mathbf{U}}^{(1)} (\boldsymbol{\Sigma}_\star^{(1)})^{-1} \right\|_{2,\infty},
\end{aligned}$$

where we have used Lemma 6 (cf. (66)) and the relation $\|\mathbf{A}\|_{2,\infty}^2 = \|\mathbf{A}\mathbf{A}^\top\|_{\infty}$ in the first equality. Using the definition of $\check{\mathbf{U}}^{(k)}$ from (4), we have

$$\begin{aligned}
\left\| \check{\mathbf{U}}^{(1)} (\boldsymbol{\Sigma}_\star^{(1)})^{-1} \right\|_{2,\infty} &= \left\| (\mathbf{U}^{(3)} \otimes \mathbf{U}^{(2)}) \mathcal{M}_1(\mathcal{G})^\top (\boldsymbol{\Sigma}_\star^{(1)})^{-1} \right\|_{2,\infty} \leq \left\| \mathbf{U}^{(3)} \otimes \mathbf{U}^{(2)} \right\|_{2,\infty} \left\| \mathcal{M}_1(\mathcal{G})^\top (\boldsymbol{\Sigma}_\star^{(1)})^{-1} \right\| \\
&\leq \left\| \mathbf{U}^{(3)} \right\|_{2,\infty} \left\| \mathbf{U}^{(2)} \right\|_{2,\infty} \left(\left\| \mathcal{M}_1(\boldsymbol{\Delta}\mathcal{G})^\top (\boldsymbol{\Sigma}_\star^{(1)})^{-1} \right\| + \left\| \mathcal{M}_1(\mathcal{G}_\star)^\top (\boldsymbol{\Sigma}_\star^{(1)})^{-1} \right\| \right).
\end{aligned}$$

Applying the triangle inequality on the decompositions $\mathbf{U}^{(k)} = \boldsymbol{\Delta}_{\mathbf{U}^{(k)}} + \mathbf{U}_\star^{(k)}$ and $\mathcal{G} = \boldsymbol{\Delta}\mathcal{G} + \mathcal{G}_\star$, and with Lemma 8 (cf. (70a)) and $\left\| \mathcal{M}_1(\mathcal{G}_\star)^\top (\boldsymbol{\Sigma}_\star^{(1)})^{-1} \right\| = 1$, it follows from the above inequalities that

$$\begin{aligned}
&\left\| \mathcal{M}_1(\mathbf{y} - \mathcal{T}_{\zeta_0}(\mathbf{y}))^\top \mathbf{U}_0^{(1)} (\mathbf{Q}_0^{(1)})^{-\top} (\boldsymbol{\Sigma}_\star^{(1)})^{-1} \right\|_{2,\infty} \\
&\leq (1 + \epsilon_0) r_1 \left(\left\| \boldsymbol{\Delta}_{\mathbf{U}^{(3)}} \right\|_{2,\infty} + \left\| \mathbf{U}_\star^{(3)} \right\|_{2,\infty} \right) \left(\left\| \boldsymbol{\Delta}_{\mathbf{U}^{(2)}} \right\|_{2,\infty} + \left\| \mathbf{U}_\star^{(2)} \right\|_{2,\infty} \right) \\
&\leq (1 + \epsilon_0) r_1 \left(\frac{\left\| \boldsymbol{\Delta}_{\mathbf{U}^{(3)}} \boldsymbol{\Sigma}_\star^{(3)} \right\|_{2,\infty}}{\sigma_{\min}(\boldsymbol{\mathcal{X}}_\star)} + \sqrt{\frac{\mu r_3}{n_3}} \right) \left(\frac{\left\| \boldsymbol{\Delta}_{\mathbf{U}^{(2)}} \boldsymbol{\Sigma}_\star^{(2)} \right\|_{2,\infty}}{\sigma_{\min}(\boldsymbol{\mathcal{X}}_\star)} + \sqrt{\frac{\mu r_2}{n_2}} \right),
\end{aligned}$$

where the last line uses the relationship $\left\| \boldsymbol{\Delta}_{\mathbf{U}^{(k)}} \right\|_{2,\infty} \leq \frac{\left\| \boldsymbol{\Delta}_{\mathbf{U}^{(k)}} \boldsymbol{\Sigma}_\star^{(k)} \right\|_{2,\infty}}{\sigma_{\min}(\boldsymbol{\Sigma}_\star^{(k)})}$ and the inequality $\left\| \mathbf{U}_\star^{(k)} \right\|_{2,\infty} \leq \sqrt{\frac{\mu r_k}{n_k}}$. Plug this back into (61) to arrive at

$$\mathfrak{A}_2 \leq \frac{0.02}{\sigma_{\min}(\boldsymbol{\mathcal{X}}_\star)} \sqrt{\frac{r_1}{\mu n_1}} \left(\sqrt{\frac{n_3}{r_3}} \left\| \boldsymbol{\Delta}_{\mathbf{U}^{(3)}} \boldsymbol{\Sigma}_\star^{(3)} \right\|_{2,\infty} + \sqrt{\mu} \sigma_{\min}(\boldsymbol{\mathcal{X}}_\star) \right) \left(\sqrt{\frac{n_2}{r_2}} \left\| \boldsymbol{\Delta}_{\mathbf{U}^{(2)}} \boldsymbol{\Sigma}_\star^{(2)} \right\|_{2,\infty} + \sqrt{\mu} \sigma_{\min}(\boldsymbol{\mathcal{X}}_\star) \right), \quad (62)$$

where we simplified the constants using the assumption $\epsilon_0 = 54.1c_0 \leq 0.15$.

Step 3: Putting things together. Combining (60), (62), and (59), we have

$$\begin{aligned}
&\frac{1}{\sigma_{\min}(\boldsymbol{\mathcal{X}}_\star)} \sqrt{\frac{n_1}{\mu r_1}} \left\| \boldsymbol{\Delta}_{\mathbf{U}^{(1)}} \boldsymbol{\Sigma}_\star^{(1)} \right\|_{2,\infty} \\
&\leq 0.57 + 0.02 \left(\frac{1}{\sigma_{\min}(\boldsymbol{\mathcal{X}}_\star)} \sqrt{\frac{n_3}{\mu r_3}} \left\| \boldsymbol{\Delta}_{\mathbf{U}^{(3)}} \boldsymbol{\Sigma}_\star^{(3)} \right\|_{2,\infty} + 1 \right) \left(\frac{1}{\sigma_{\min}(\boldsymbol{\mathcal{X}}_\star)} \sqrt{\frac{n_2}{\mu r_2}} \left\| \boldsymbol{\Delta}_{\mathbf{U}^{(2)}} \boldsymbol{\Sigma}_\star^{(2)} \right\|_{2,\infty} + 1 \right).
\end{aligned}$$

Similar inequalities hold for $\sqrt{\frac{n_2}{r_2}} \left\| \boldsymbol{\Delta}_{\mathbf{U}^{(2)}} \boldsymbol{\Sigma}_\star^{(2)} \right\|_{2,\infty}$ and $\sqrt{\frac{n_3}{r_3}} \left\| \boldsymbol{\Delta}_{\mathbf{U}^{(3)}} \boldsymbol{\Sigma}_\star^{(3)} \right\|_{2,\infty}$. Taking the maximum as $\mathcal{I} := \max_k \frac{1}{\sigma_{\min}(\boldsymbol{\mathcal{X}}_\star)} \sqrt{\frac{n_k}{\mu r_k}} \left\| \boldsymbol{\Delta}_{\mathbf{U}^{(k)}} \boldsymbol{\Sigma}_\star^{(k)} \right\|_{2,\infty}$, we have

$$\mathcal{I} \leq 0.57 + 0.02(\mathcal{I} + 1)^2 \quad \implies \quad \mathcal{I} \leq 0.62,$$

and consequently, $\max_k \sqrt{\frac{n_k}{r_k}} \left\| \boldsymbol{\Delta}_{\mathbf{U}^{(k)}} \boldsymbol{\Sigma}_\star^{(k)} \right\|_{2,\infty} < \sqrt{\mu} \sigma_{\min}(\boldsymbol{\mathcal{X}}_\star)$ as claimed.

We are left with proving the decomposition (58).

Proof of (58) . By the assumption

$$\alpha \leq \frac{c_0}{\mu^2 r_1 r_2 r_3 \kappa} \leq \frac{c_0}{\sqrt{\mu^3 r_1 r_2 r_3 r \kappa}},$$

in view of Lemma 3, we have that

$$\text{dist}(\mathbf{F}_0, \mathbf{F}_\star) < 54.1 c_0 \sigma_{\min}(\mathcal{X}_\star) =: \epsilon_0 \sigma_{\min}(\mathcal{X}_\star) \leq 0.15 \sigma_{\min}(\mathcal{X}_\star), \quad (63)$$

where $\epsilon_0 = 54.1 c_0 \leq 0.15$ as long as $c_0 > 0$ is small enough. Given $\text{dist}(\mathbf{F}_0, \mathbf{F}_\star) < \sigma_{\min}(\mathcal{X}_\star)$, we know that $\{\mathbf{Q}_0^{(k)}\}_{k=1}^3$, the optimal alignment matrices between \mathbf{F}_0 and \mathbf{F}_\star exist by [TMPB⁺22, Lemma 6].

We now aim to control the incoherence. We begin with the equality guaranteed by Lemma 6 (cf. (66)),

$$\mathbf{U}_0^{(k)} \check{\mathbf{U}}_0^{(k)\top} \check{\mathbf{U}}_0^{(k)} = \mathcal{M}_k(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y})) \mathcal{M}_k(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y}))^\top \mathbf{U}_0^{(k)}.$$

Again, we will focus on the case with $k = 1$; the other modes will follow from the same arguments. Given that $(\check{\mathbf{U}}_0^{(1)\top} \check{\mathbf{U}}_0^{(1)})^{-1}$ exists since $\check{\mathbf{U}}_0^{(1)\top} \check{\mathbf{U}}_0^{(1)} = \mathcal{M}_1(\mathcal{G}_0) \mathcal{M}_1(\mathcal{G}_0)^\top$ is positive definite, right-multiplying $(\check{\mathbf{U}}_0^{(1)\top} \check{\mathbf{U}}_0^{(1)})^{-1}$ on both sides of the above equation yields

$$\mathbf{U}_0^{(1)} = \mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y})) \mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y}))^\top \mathbf{U}_0^{(1)} (\check{\mathbf{U}}_0^{(1)\top} \check{\mathbf{U}}_0^{(1)})^{-1}.$$

Plug in the relation $\mathcal{Y} = \mathcal{X}_\star + \mathcal{S}_\star$ to get

$$\begin{aligned} \mathbf{U}_0^{(1)} &= \mathcal{M}_1(\mathcal{S}_\star - \mathcal{T}_{\zeta_0}(\mathcal{Y})) \mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y}))^\top \mathbf{U}_0^{(1)} (\check{\mathbf{U}}_0^{(1)\top} \check{\mathbf{U}}_0^{(1)})^{-1} \\ &\quad + \mathcal{M}_1(\mathcal{X}_\star) \mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y}))^\top \mathbf{U}_0^{(1)} (\check{\mathbf{U}}_0^{(1)\top} \check{\mathbf{U}}_0^{(1)})^{-1} \\ &= \mathcal{M}_1(\mathcal{S}_\star - \mathcal{T}_{\zeta_0}(\mathcal{Y})) \mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y}))^\top \mathbf{U}_0^{(1)} (\check{\mathbf{U}}_0^{(1)\top} \check{\mathbf{U}}_0^{(1)})^{-1} \\ &\quad + \mathbf{U}_\star^{(1)} \check{\mathbf{U}}_\star^{(1)\top} \mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y}))^\top \mathbf{U}_0^{(1)} (\check{\mathbf{U}}_0^{(1)\top} \check{\mathbf{U}}_0^{(1)})^{-1}. \end{aligned}$$

Subtracting $\mathbf{U}_\star^{(1)} (\mathbf{Q}_0^{(1)})^{-1}$ on both sides gets us

$$\begin{aligned} \mathbf{U}_0^{(1)} - \mathbf{U}_\star^{(1)} (\mathbf{Q}_0^{(1)})^{-1} &= \mathcal{M}_1(\mathcal{S}_\star - \mathcal{T}_{\zeta_0}(\mathcal{Y})) \mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y}))^\top \mathbf{U}_0^{(1)} (\check{\mathbf{U}}_0^{(1)\top} \check{\mathbf{U}}_0^{(1)})^{-1} \\ &\quad + \mathbf{U}_\star^{(1)} \check{\mathbf{U}}_\star^{(1)\top} \mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y}))^\top \mathbf{U}_0^{(1)} (\check{\mathbf{U}}_0^{(1)\top} \check{\mathbf{U}}_0^{(1)})^{-1} - \mathbf{U}_\star^{(1)} (\mathbf{Q}_0^{(1)})^{-1}. \end{aligned} \quad (64)$$

Observe that

$$\begin{aligned} \mathbf{U}_\star^{(1)} (\mathbf{Q}_0^{(1)})^{-1} &= \mathbf{U}_\star^{(1)} (\mathbf{Q}_0^{(1)})^{-1} \check{\mathbf{U}}_0^{(1)\top} \check{\mathbf{U}}_0^{(1)} (\check{\mathbf{U}}_0^{(1)\top} \check{\mathbf{U}}_0^{(1)})^{-1} \\ &= \mathbf{U}_\star^{(1)} (\mathbf{Q}_0^{(1)})^{-1} \check{\mathbf{U}}_0^{(1)\top} \mathcal{M}_1(\mathcal{Y} - \mathcal{T}_{\zeta_0}(\mathcal{Y}))^\top \mathbf{U}_0^{(1)} (\check{\mathbf{U}}_0^{(1)\top} \check{\mathbf{U}}_0^{(1)})^{-1}, \end{aligned}$$

where we have used Lemma 6 (cf. (65)) in the second step. Plug this result into (64), multiply both sides with $\mathbf{Q}_0^{(1)} \Sigma_\star^{(1)}$, and recall the short-hand notation in (21) initiated at $t = 0$ to arrive at the claimed decomposition.

E Technical lemmas

This section collects several technical lemmas that are useful in the main proofs.

E.1 Tensor algebra

We start with a simple bound on the element-wise maximum norm of an incoherent tensor.

Lemma 5. *Suppose that $\mathcal{X}_\star = (\mathbf{U}_\star^{(1)}, \mathbf{U}_\star^{(2)}, \mathbf{U}_\star^{(3)}) \cdot \mathcal{G}_\star \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ have multilinear rank $\mathbf{r} = (r_1, r_2, r_3)$ and is μ -incoherent. Then one has $\|\mathcal{X}_\star\|_\infty \leq \sqrt{\frac{\mu^3 r_1 r_2 r_3}{n_1 n_2 n_3}} \kappa \sigma_{\min}(\mathcal{X}_\star)$.*

Proof. By the property of matricizations (cf. (4)), we have $\mathcal{M}_k(\mathcal{X}_*) = \mathbf{U}_*^{(k)} \check{\mathbf{U}}_*^{(k)\top}$ for any $k = 1, 2, 3$. Furthermore, $\|\cdot\|_\infty$ is invariant to matricizations, so

$$\|\mathcal{X}_*\|_\infty = \|\mathcal{M}_k(\mathcal{X}_*)\|_\infty = \left\| \mathbf{U}_*^{(k)} \check{\mathbf{U}}_*^{(k)\top} \right\|_\infty.$$

Without loss of generality, we choose $k = 1$. It then follows that

$$\begin{aligned} \|\mathcal{X}_*\|_\infty &\leq \left\| \mathbf{U}_*^{(1)} \right\|_{2,\infty} \left\| \check{\mathbf{U}}_*^{(1)} \right\|_{2,\infty} \leq \left\| \mathbf{U}_*^{(1)} \right\|_{2,\infty} \left\| \mathbf{U}_*^{(3)} \right\|_{2,\infty} \left\| \mathbf{U}_*^{(2)} \right\|_{2,\infty} \|\mathcal{M}_1(\mathcal{G}_*)\| \\ &\leq \sqrt{\frac{\mu^3 r_1 r_2 r_3}{n_1 n_2 n_3}} \sigma_{\max}(\mathcal{M}_1(\mathcal{X}_*)), \end{aligned}$$

where the second line follows from the definition of incoherence of \mathcal{X}_* , i.e., $\left\| \mathbf{U}_*^{(k)} \right\|_{2,\infty} \leq \sqrt{\frac{\mu r_k}{n_k}}$, and $\|\mathcal{M}_1(\mathcal{G}_*)\| = \sigma_{\max}(\mathcal{M}_1(\mathcal{X}_*))$. Applying the above bound for any k and taking the tightest bound, we have

$$\|\mathcal{X}_*\|_\infty \leq \sqrt{\frac{\mu^3 r_1 r_2 r_3}{n_1 n_2 n_3}} \min_k \sigma_{\max}(\mathcal{M}_k(\mathcal{X}_*)) = \sqrt{\frac{\mu^3 r_1 r_2 r_3}{n_1 n_2 n_3}} \kappa \sigma_{\min}(\mathcal{X}_*),$$

where we have used the definition of κ . \square

We next show a key tensor algebraic result that is crucial in establishing the incoherence property of the spectral initialization.

Lemma 6. *Given a tensor $\mathcal{T} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$, suppose its rank- \mathbf{r} truncated HOSVD is $(\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathcal{G}$ with $\mathbf{U}^{(k)} \in \mathbb{R}^{n_k \times r_k}$ and $\mathcal{G} \in \mathbb{R}^{r_1 \times r_2 \times r_3}$, and $r_k \leq n_k$ for $k = 1, 2, 3$. Then,*

$$\mathbf{U}^{(k)\top} \mathcal{M}_k(\mathcal{T}) \check{\mathbf{U}}^{(k)} = \check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)}, \quad (65)$$

$$\mathcal{M}_k(\mathcal{T}) \mathcal{M}_k(\mathcal{T})^\top \mathbf{U}^{(k)} = \mathbf{U}^{(k)} \check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)}, \quad (66)$$

where $\check{\mathbf{U}}^{(k)}$ is defined in (4).

Proof. Set the full HOSVD of \mathcal{T} to be $(\mathbf{U}_\mathcal{T}^{(1)}, \mathbf{U}_\mathcal{T}^{(2)}, \mathbf{U}_\mathcal{T}^{(3)}) \cdot \mathcal{G}_\mathcal{T}$. Since $\mathbf{U}^{(k)}$ contains all the left singular vectors of $\mathcal{M}_k(\mathcal{T})$, it can be decomposed into the following block structure:

$$\mathbf{U}^{(k)} = [\mathbf{U}^{(k)} \quad \bar{\mathbf{U}}^{(k)}], \quad (67)$$

where $\bar{\mathbf{U}}^{(k)} \in \mathbb{R}^{n_k \times (n_k - r_k)}$ contains the bottom $(n_k - r_k)$ left singular vectors. The rest of the proof focuses on the first mode (i.e., $k = 1$), while other modes follow from similar arguments.

Let us begin with proving (65). Plugging in the definition of $\check{\mathbf{U}}^{(1)}$ from (4a), we see that

$$\begin{aligned} \mathbf{U}^{(1)\top} \mathcal{M}_1(\mathcal{T}) \check{\mathbf{U}}^{(1)} &= \mathbf{U}^{(1)\top} \mathbf{U}_\mathcal{T}^{(1)} \mathcal{M}_1(\mathcal{G}_\mathcal{T}) (\mathbf{U}_\mathcal{T}^{(3)} \otimes \mathbf{U}_\mathcal{T}^{(2)})^\top (\mathbf{U}^{(3)} \otimes \mathbf{U}^{(2)}) \mathcal{M}_1(\mathcal{G})^\top \\ &= \mathbf{U}^{(1)\top} [\mathbf{U}^{(1)} \quad \bar{\mathbf{U}}^{(1)}] \mathcal{M}_1(\mathcal{G}_\mathcal{T}) \left(\begin{bmatrix} \mathbf{U}^{(3)\top} \\ \bar{\mathbf{U}}^{(3)\top} \end{bmatrix} \otimes \begin{bmatrix} \mathbf{U}^{(2)\top} \\ \bar{\mathbf{U}}^{(2)\top} \end{bmatrix} \right) (\mathbf{U}^{(3)} \otimes \mathbf{U}^{(2)}) \mathcal{M}_1(\mathcal{G})^\top, \end{aligned}$$

where the second line uses the block structure (67). By the mixed product property of Kronecker products, we have

$$\begin{aligned} \mathbf{U}^{(1)\top} \mathcal{M}_1(\mathcal{T}) \check{\mathbf{U}}^{(1)} &= [\mathbf{I}_{r_1} \quad \mathbf{0}] \mathcal{M}_1(\mathcal{G}_\mathcal{T}) \left(\begin{bmatrix} \mathbf{U}^{(3)\top} \mathbf{U}^{(3)} \\ \bar{\mathbf{U}}^{(3)\top} \mathbf{U}^{(3)} \end{bmatrix} \otimes \begin{bmatrix} \mathbf{U}^{(2)\top} \mathbf{U}^{(2)} \\ \bar{\mathbf{U}}^{(2)\top} \mathbf{U}^{(2)} \end{bmatrix} \right) \mathcal{M}_1(\mathcal{G})^\top \\ &= [\mathbf{I}_{r_1} \quad \mathbf{0}] \mathcal{M}_1(\mathcal{G}_\mathcal{T}) \left(\begin{bmatrix} \mathbf{I}_{r_3} \\ \mathbf{0} \end{bmatrix} \otimes \begin{bmatrix} \mathbf{I}_{r_2} \\ \mathbf{0} \end{bmatrix} \right) \mathcal{M}_1(\mathcal{G})^\top, \end{aligned} \quad (68)$$

where we used the fact that the singular vectors are orthonormal. Note that

$$[\mathbf{I}_{r_1} \ \mathbf{0}] \mathcal{M}_1(\mathcal{G}_{\mathcal{T}}) \left(\begin{bmatrix} \mathbf{I}_{r_3} \\ \mathbf{0} \end{bmatrix} \otimes \begin{bmatrix} \mathbf{I}_{r_2} \\ \mathbf{0} \end{bmatrix} \right) = \mathcal{M}_1 \left(([\mathbf{I}_{r_1} \ \mathbf{0}], [\mathbf{I}_{r_2} \ \mathbf{0}], [\mathbf{I}_{r_3} \ \mathbf{0}]) \cdot \mathcal{G}_{\mathcal{T}} \right) \in \mathbb{R}^{r_1 \times r_2 r_3},$$

where $([\mathbf{I}_{r_1} \ \mathbf{0}], [\mathbf{I}_{r_2} \ \mathbf{0}], [\mathbf{I}_{r_3} \ \mathbf{0}]) \cdot \mathcal{G}_{\mathcal{T}}$ is equivalent to trimming off the entries $[\mathcal{G}_{\mathcal{T}}]_{i_1, i_2, i_3}$ for all $i_1 > r_1$, $i_2 > r_2$, or $i_3 > r_3$. Since \mathcal{G} is the section of $[\mathcal{G}_{\mathcal{T}}]_{i_1, i_2, i_3}$ where $1 \leq i_1 \leq r_1, 1 \leq i_2 \leq r_2, 1 \leq i_3 \leq r_3$ [VVM12], we have

$$([\mathbf{I}_{r_1} \ \mathbf{0}], [\mathbf{I}_{r_2} \ \mathbf{0}], [\mathbf{I}_{r_3} \ \mathbf{0}]) \cdot \mathcal{G}_{\mathcal{T}} = \mathcal{G}.$$

This allows us to simplify (68) as

$$\mathbf{U}^{(1)\top} \mathcal{M}_1(\mathcal{T}) \check{\mathbf{U}}^{(1)} = \mathcal{M}_1(\mathcal{G}) \mathcal{M}_1(\mathcal{G})^\top = \check{\mathbf{U}}^{(1)\top} \check{\mathbf{U}}^{(1)},$$

where the last equality follows from the definition of $\check{\mathbf{U}}^{(1)}$ from (4) and $(\mathbf{U}^{(3)} \otimes \mathbf{U}^{(2)})^\top (\mathbf{U}^{(3)} \otimes \mathbf{U}^{(2)}) = \mathbf{I}$ by construction of the HOSVD. This completes the proof of (65).

Turning to (66), we begin with the observation

$$\begin{aligned} \mathcal{M}_1(\mathcal{T}) \mathcal{M}_1(\mathcal{T})^\top &= \mathbf{U}_{\mathcal{T}}^{(1)} \mathcal{M}_1(\mathcal{G}_{\mathcal{T}}) (\mathbf{U}_{\mathcal{T}}^{(3)} \otimes \mathbf{U}_{\mathcal{T}}^{(2)})^\top (\mathbf{U}_{\mathcal{T}}^{(3)} \otimes \mathbf{U}_{\mathcal{T}}^{(2)}) \mathcal{M}_1(\mathcal{G}_{\mathcal{T}})^\top \mathbf{U}_{\mathcal{T}}^{(1)\top} \\ &= \mathbf{U}_{\mathcal{T}}^{(1)} \mathcal{M}_1(\mathcal{G}_{\mathcal{T}}) \mathcal{M}_1(\mathcal{G}_{\mathcal{T}})^\top \mathbf{U}_{\mathcal{T}}^{(1)\top}. \end{aligned} \quad (69)$$

By the ‘‘all-orthogonal’’ property in [DLDMV00], $\mathcal{M}_1(\mathcal{G}_{\mathcal{T}}) \mathcal{M}_1(\mathcal{G}_{\mathcal{T}})^\top = (\boldsymbol{\Sigma}_{\mathcal{T}}^{(1)})^2$, the squared singular value matrix of $\mathcal{M}_1(\mathcal{T})$. By assigning $\boldsymbol{\Sigma}_{\mathcal{T}}^{(1)}$ with the block representation

$$\boldsymbol{\Sigma}_{\mathcal{T}}^{(1)} = \begin{bmatrix} \boldsymbol{\Sigma}^{(1)} & \mathbf{0} \\ \mathbf{0} & \bar{\boldsymbol{\Sigma}}^{(1)} \end{bmatrix},$$

where $\boldsymbol{\Sigma}^{(1)} = \sqrt{\mathcal{M}_1(\mathcal{G}) \mathcal{M}_1(\mathcal{G})^\top}$ and $\bar{\boldsymbol{\Sigma}}^{(1)}$ contain the top r_1 singular values and bottom $n_1 - r_1$ singular values, respectively, of $\mathcal{M}_1(\mathcal{T})$. Coupled with same block structure as in (67), (69) becomes

$$\mathcal{M}_1(\mathcal{T}) \mathcal{M}_1(\mathcal{T})^\top = [\mathbf{U}^{(1)} \ \bar{\mathbf{U}}^{(1)}] \begin{bmatrix} (\boldsymbol{\Sigma}^{(1)})^2 & \mathbf{0} \\ \mathbf{0} & (\bar{\boldsymbol{\Sigma}}^{(1)})^2 \end{bmatrix} \begin{bmatrix} \mathbf{U}^{(1)\top} \\ \bar{\mathbf{U}}^{(1)\top} \end{bmatrix}.$$

Multiply by $\mathbf{U}^{(1)}$ on the right to arrive at $\mathcal{M}_1(\mathcal{T}) \mathcal{M}_1(\mathcal{T})^\top \mathbf{U}^{(1)} = \mathbf{U}^{(1)} (\boldsymbol{\Sigma}^{(1)})^2 = \mathbf{U}^{(1)} \check{\mathbf{U}}^{(1)\top} \check{\mathbf{U}}^{(1)}$. \square

E.2 Perturbation bounds

Below is a useful perturbation bound for matrices.

Lemma 7. *Given two matrices $\mathbf{U}, \mathbf{U}_* \in \mathbb{R}^{n \times r}$ that have full column rank, two invertible matrices $\bar{\mathbf{Q}}, \mathbf{Q} \in \mathbb{R}^{r \times r}$, and a positive definite matrix $\boldsymbol{\Sigma} \in \mathbb{R}^{r \times r}$. Suppose that $\sigma_{\min}(\mathbf{U}_*) > \|\mathbf{U}\bar{\mathbf{Q}} - \mathbf{U}_*\|$. Then the following holds true*

$$\|\bar{\mathbf{Q}}^{-1}(\mathbf{Q} - \bar{\mathbf{Q}})\boldsymbol{\Sigma}\| \leq \frac{\|\mathbf{U}(\mathbf{Q} - \bar{\mathbf{Q}})\boldsymbol{\Sigma}\|}{\sigma_{\min}(\mathbf{U}_*) - \|\mathbf{U}\bar{\mathbf{Q}} - \mathbf{U}_*\|}.$$

Proof. It follows that

$$\begin{aligned} \|\bar{\mathbf{Q}}^{-1}(\mathbf{Q} - \bar{\mathbf{Q}})\boldsymbol{\Sigma}\| &= \|\bar{\mathbf{Q}}^{-1}(\mathbf{U}^\top \mathbf{U})^{-1} \mathbf{U}^\top \mathbf{U}(\mathbf{Q} - \bar{\mathbf{Q}})\boldsymbol{\Sigma}\| \\ &\leq \|\bar{\mathbf{Q}}^{-1}(\mathbf{U}^\top \mathbf{U})^{-1} \mathbf{U}^\top\| \|\mathbf{U}(\mathbf{Q} - \bar{\mathbf{Q}})\boldsymbol{\Sigma}\| \\ &= \frac{\|\mathbf{U}(\mathbf{Q} - \bar{\mathbf{Q}})\boldsymbol{\Sigma}\|}{\sigma_{\min}(\mathbf{U}\bar{\mathbf{Q}})}, \end{aligned}$$

where the last equality comes from the fact that $\bar{Q}^{-1}(\mathbf{U}^\top \mathbf{U})^{-1} \mathbf{U}^\top$ and $\mathbf{U} \bar{Q}$ are pseudoinverses. By Weyl's inequality, we know $|\sigma_i(\mathbf{M}) - \sigma_i(\mathbf{M} + \Delta_{\mathbf{M}})| \leq \|\Delta_{\mathbf{M}}\|$. Taking $\mathbf{M} := \mathbf{U}_*$ and $\Delta_{\mathbf{M}} := \mathbf{U} \bar{Q} - \mathbf{U}_*$, we have

$$\|\bar{Q}^{-1}(\mathbf{Q} - \bar{Q})\Sigma\| \leq \frac{\|\mathbf{U}(\mathbf{Q} - \bar{Q})\Sigma\|}{\sigma_{\min}(\mathbf{U}_*) - \|\mathbf{U} \bar{Q} - \mathbf{U}_*\|}$$

as long as $\sigma_{\min}(\mathbf{U}_*) > \|\mathbf{U} \bar{Q} - \mathbf{U}_*\|$. \square

We also collect a useful lemma regarding perturbation bounds for tensors from [TMPB+22].

Lemma 8 ([TMPB+22, Lemma 10]). *Suppose $\mathbf{F} = (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}, \mathcal{G})$ and $\mathbf{F}_* = (\mathbf{U}_*^{(1)}, \mathbf{U}_*^{(2)}, \mathbf{U}_*^{(3)}, \mathcal{G}_*)$ are aligned, and $\text{dist}(\mathbf{F}, \mathbf{F}_*) \leq \epsilon \sigma_{\min}(\mathcal{X}_*)$ for some $0 < \epsilon < 1$. Then, the following bounds are true:*

$$\|\mathcal{M}_k(\Delta_{\mathcal{G}})^\top (\Sigma_*^{(k)})^{-1}\| \leq \epsilon; \quad (70a)$$

$$\|\mathbf{U}^{(k)}(\mathbf{U}^{(k)\top} \mathbf{U}^{(k)})^{-1}\| \leq \frac{1}{1 - \epsilon}; \quad (70b)$$

$$\|\check{\mathbf{U}}^{(k)}(\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} \Sigma_*^{(k)}\| \leq \frac{1}{(1 - \epsilon)^3}; \quad (70c)$$

$$\|\Sigma_*^{(k)}(\check{\mathbf{U}}^{(k)\top} \check{\mathbf{U}}^{(k)})^{-1} \Sigma_*^{(k)}\| \leq \frac{1}{(1 - \epsilon)^6}; \quad (70d)$$

$$\|\check{\mathbf{U}}^{(1)} - \check{\mathbf{U}}_*^{(1)}\|_{\text{F}} \leq \left(1 + \epsilon + \frac{\epsilon^2}{3}\right) \left(\|(\mathbf{U}^{(2)} - \mathbf{U}_*^{(2)})\Sigma_*^{(2)}\|_{\text{F}} + \|(\mathbf{U}^{(3)} - \mathbf{U}_*^{(3)})\Sigma_*^{(3)}\|_{\text{F}} + \|\mathcal{G} - \mathcal{G}_*\|_{\text{F}}\right). \quad (70e)$$

For (70e), similar bounds exist for the other modes. Furthermore, if $0 < \epsilon \leq 0.2$,

$$\|(\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathcal{G} - \mathcal{X}_*\|_{\text{F}} \leq 3 \text{dist}(\mathbf{F}, \mathbf{F}_*). \quad (71)$$

The next set of lemmas, which is crucial in our analysis, deals with perturbation bounds when relating a tensor $\mathcal{X} = (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathcal{G}$ to the ground truth \mathcal{X}_* , where the tensor tuples $\mathbf{F} = (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}, \mathcal{G})$ and $\mathbf{F}_* = (\mathbf{U}_*^{(1)}, \mathbf{U}_*^{(2)}, \mathbf{U}_*^{(3)}, \mathcal{G}_*)$ are aligned.

Lemma 9. *Let $\mathcal{X}_* \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ be μ -incoherent with the Tucker decomposition $\mathcal{X}_* = (\mathbf{U}_*^{(1)}, \mathbf{U}_*^{(2)}, \mathbf{U}_*^{(3)}) \cdot \mathcal{G}_*$ of rank $\mathbf{r} = (r_1, r_2, r_3)$, and $\{\Sigma_*^{(k)}\}_{k=1,2,3}$ be the set of singular value matrices of different matricizations of \mathcal{X}_* . In addition, let $\mathbf{F} := (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}, \mathcal{G})$ and $\mathbf{F}_* := (\mathbf{U}_*^{(1)}, \mathbf{U}_*^{(2)}, \mathbf{U}_*^{(3)}, \mathcal{G}_*)$ be aligned, where $\mathcal{X} = (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathcal{G}$. Suppose*

$$\max_k \left\{ \sqrt{\frac{n_k}{r_k}} \left\| (\mathbf{U}^{(k)} - \mathbf{U}_*^{(k)}) \Sigma_*^{(k)} \right\|_{2,\infty} \right\} \leq c \sqrt{\mu} \sigma_{\min}(\mathcal{X}_*) \quad (72)$$

for some $0 < c \leq 1$. Then for $k = 1, 2, 3$,

$$\|\mathbf{U}^{(k)} - \mathbf{U}_*^{(k)}\|_{2,\infty} \leq c \sqrt{\frac{\mu r_k}{n_k}}, \quad \text{and} \quad \|\mathbf{U}^{(k)}\|_{2,\infty} \leq 2 \sqrt{\frac{\mu r_k}{n_k}}. \quad (73)$$

Proof. It follows that for all k ,

$$\|\mathbf{U}^{(k)} - \mathbf{U}_*^{(k)}\|_{2,\infty} \leq \frac{1}{\sigma_{\min}(\Sigma_*^{(k)})} \left\| (\mathbf{U}^{(k)} - \mathbf{U}_*^{(k)}) \Sigma_*^{(k)} \right\|_{2,\infty} \leq c \sqrt{\frac{\mu r_k}{n_k}},$$

where the second inequality follows from (72) and $\sigma_{\min}(\mathcal{X}_*) \leq \sigma_{\min}(\Sigma_*^{(k)})$. This completes the proof for the first part of (73). With this and the incoherence assumption $\|\mathbf{U}_*^{(k)}\|_{2,\infty} \leq \sqrt{\frac{\mu r_k}{n_k}}$, after applying triangle inequality, we arrive at

$$\|\mathbf{U}^{(k)}\|_{2,\infty} \leq \|\mathbf{U}^{(k)} - \mathbf{U}_*^{(k)}\|_{2,\infty} + \|\mathbf{U}_*^{(k)}\|_{2,\infty} \leq 2 \sqrt{\frac{\mu r_k}{n_k}},$$

which completes the proof. \square

Lemma 10. Let $\mathcal{X}_* \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ be μ -incoherent with the Tucker decomposition $\mathcal{X}_* = (\mathbf{U}_*^{(1)}, \mathbf{U}_*^{(2)}, \mathbf{U}_*^{(3)}) \cdot \mathcal{G}_*$ of rank $\mathbf{r} = (r_1, r_2, r_3)$. In addition, let $\mathbf{F} := (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}, \mathcal{G})$ and $\mathbf{F}_* := (\mathbf{U}_*^{(1)}, \mathbf{U}_*^{(2)}, \mathbf{U}_*^{(3)}, \mathcal{G}_*)$ be aligned, where $\mathcal{X} = (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathcal{G}$. For $0 < \epsilon_0 < 0.1$ and $0 < c \leq 1$, if

$$\text{dist}(\mathbf{F}, \mathbf{F}_*) \leq \epsilon_0 c \sigma_{\min}(\mathcal{X}_*), \quad (74a)$$

$$\max_k \left\{ \sqrt{\frac{n_k}{r_k}} \left\| (\mathbf{U}^{(k)} - \mathbf{U}_*^{(k)}) \boldsymbol{\Sigma}_*^{(k)} \right\|_{2, \infty} \right\} \leq c \sqrt{\mu} \sigma_{\min}(\mathcal{X}_*) \quad (74b)$$

are satisfied, then

$$\|\mathcal{X} - \mathcal{X}_*\|_{\infty} \leq \sqrt{\frac{\mu^3 r_1 r_2 r_3}{n_1 n_2 n_3}} (8\epsilon_0 + 7) c \sigma_{\min}(\mathcal{X}_*) \leq 8 \sqrt{\frac{\mu^3 r_1 r_2 r_3}{n_1 n_2 n_3}} c \sigma_{\min}(\mathcal{X}_*).$$

Proof. We can decompose $\mathcal{X} - \mathcal{X}_*$ into

$$\begin{aligned} \mathcal{X} - \mathcal{X}_* &= (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot (\mathcal{G} - \mathcal{G}_*) + (\mathbf{U}^{(1)} - \mathbf{U}_*^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathcal{G}_* \\ &\quad + (\mathbf{U}_*^{(1)}, \mathbf{U}^{(2)} - \mathbf{U}_*^{(2)}, \mathbf{U}^{(3)}) \cdot \mathcal{G}_* + (\mathbf{U}_*^{(1)}, \mathbf{U}_*^{(2)}, \mathbf{U}^{(3)} - \mathbf{U}_*^{(3)}) \cdot \mathcal{G}_*. \end{aligned} \quad (75)$$

Then by the triangle inequality,

$$\begin{aligned} \|\mathcal{X} - \mathcal{X}_*\|_{\infty} &\leq \left\| (\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot (\mathcal{G} - \mathcal{G}_*) \right\|_{\infty} + \left\| (\mathbf{U}^{(1)} - \mathbf{U}_*^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) \cdot \mathcal{G}_* \right\|_{\infty} \\ &\quad + \left\| (\mathbf{U}_*^{(1)}, \mathbf{U}^{(2)} - \mathbf{U}_*^{(2)}, \mathbf{U}^{(3)}) \cdot \mathcal{G}_* \right\|_{\infty} + \left\| (\mathbf{U}_*^{(1)}, \mathbf{U}_*^{(2)}, \mathbf{U}^{(3)} - \mathbf{U}_*^{(3)}) \cdot \mathcal{G}_* \right\|_{\infty} \\ &= \underbrace{\left\| \mathbf{U}^{(1)} \mathcal{M}_1(\mathcal{G} - \mathcal{G}_*) (\mathbf{U}^{(3)} \otimes \mathbf{U}^{(2)})^{\top} \right\|_{\infty}}_{=: \mathfrak{A}_{\text{core}}} + \underbrace{\left\| (\mathbf{U}^{(1)} - \mathbf{U}_*^{(1)}) \mathcal{M}_1(\mathcal{G}_*) (\mathbf{U}^{(3)} \otimes \mathbf{U}^{(2)})^{\top} \right\|_{\infty}}_{=: \mathfrak{A}_1} \\ &\quad + \underbrace{\left\| (\mathbf{U}^{(2)} - \mathbf{U}_*^{(2)}) \mathcal{M}_2(\mathcal{G}_*) (\mathbf{U}^{(3)} \otimes \mathbf{U}_*^{(1)})^{\top} \right\|_{\infty}}_{=: \mathfrak{A}_2} + \underbrace{\left\| (\mathbf{U}^{(3)} - \mathbf{U}_*^{(3)}) \mathcal{M}_3(\mathcal{G}_*) (\mathbf{U}_*^{(2)} \otimes \mathbf{U}_*^{(1)})^{\top} \right\|_{\infty}}_{=: \mathfrak{A}_3}, \end{aligned}$$

where the second inequality follows from the invariance of ℓ_{∞} norm to matricizations. We will bound each term separately.

- For $\mathfrak{A}_{\text{core}}$, it follows from basic norm relations that

$$\begin{aligned} \mathfrak{A}_{\text{core}} &\leq \left\| \mathbf{U}^{(1)} \right\|_{2, \infty} \left\| \mathcal{M}_1(\mathcal{G} - \mathcal{G}_*) \right\| \left\| \mathbf{U}^{(3)} \otimes \mathbf{U}^{(2)} \right\|_{2, \infty} \\ &\leq \left\| \mathbf{U}^{(1)} \right\|_{2, \infty} \left\| \mathbf{U}^{(2)} \right\|_{2, \infty} \left\| \mathbf{U}^{(3)} \right\|_{2, \infty} \|\mathcal{G} - \mathcal{G}_*\|_{\text{F}} \leq 8 \sqrt{\frac{\mu^3 r_1 r_2 r_3}{n_1 n_2 n_3}} \epsilon_0 c \sigma_{\min}(\mathcal{X}_*), \end{aligned}$$

where the last inequality follows from $\|\mathcal{G} - \mathcal{G}_*\|_{\text{F}} \leq \text{dist}(\mathbf{F}, \mathbf{F}_*) \leq \epsilon_0 c \sigma_{\min}(\mathcal{X}_*)$ by assumption (74a) and Lemma 9 by assumption (74b).

- Next, for \mathfrak{A}_1 ,

$$\begin{aligned} \mathfrak{A}_1 &\leq \left\| (\mathbf{U}^{(1)} - \mathbf{U}_*^{(1)}) \boldsymbol{\Sigma}_*^{(1)} \right\|_{2, \infty} \left\| \mathbf{U}^{(2)} \right\|_{2, \infty} \left\| \mathbf{U}^{(3)} \right\|_{2, \infty} \left\| \mathcal{M}_1(\mathcal{G}_*)^{\top} (\boldsymbol{\Sigma}_*^{(1)})^{-1} \right\| \\ &= \left\| (\mathbf{U}^{(1)} - \mathbf{U}_*^{(1)}) \boldsymbol{\Sigma}_*^{(1)} \right\|_{2, \infty} \left\| \mathbf{U}^{(2)} \right\|_{2, \infty} \left\| \mathbf{U}^{(3)} \right\|_{2, \infty} \leq 4 \sqrt{\frac{\mu^3 r_1 r_2 r_3}{n_1 n_2 n_3}} c \sigma_{\min}(\mathcal{X}_*), \end{aligned}$$

where the equality follows from $\left\| \mathcal{M}_k(\mathcal{G}_*)^{\top} (\boldsymbol{\Sigma}_*^{(k)})^{-1} \right\| = 1$ since $\mathcal{M}_k(\mathcal{G}_*) \mathcal{M}_k(\mathcal{G}_*)^{\top} = (\boldsymbol{\Sigma}_*^{(k)})^2$, and the last inequality follows from the assumption (74b) and Lemma 9 by assumption (74b).

- Similarly, for \mathfrak{A}_2 , it follows

$$\begin{aligned}\mathfrak{A}_2 &\leq \left\| (U^{(2)} - U_\star^{(2)}) \Sigma_\star^{(2)} \right\|_{2,\infty} \left\| U^{(3)} \right\|_{2,\infty} \left\| U_\star^{(1)} \right\|_{2,\infty} \left\| \mathcal{M}_2(\mathcal{G}_\star)^\top (\Sigma_\star^{(2)})^{-1} \right\| \\ &\leq 2 \sqrt{\frac{\mu^3 r_1 r_2 r_3}{n_1 n_2 n_3}} c \sigma_{\min}(\mathcal{X}_\star).\end{aligned}$$

- Finally, repeat the same approach for \mathfrak{A}_3 to get

$$\begin{aligned}\mathfrak{A}_3 &\leq \left\| (U^{(3)} - U_\star^{(3)}) \Sigma_\star^{(3)} \right\|_{2,\infty} \left\| U_\star^{(2)} \right\|_{2,\infty} \left\| U_\star^{(1)} \right\|_{2,\infty} \left\| \mathcal{M}_3(\mathcal{G}_\star)^\top (\Sigma_\star^{(3)})^{-1} \right\| \\ &\leq \sqrt{\frac{\mu^3 r_1 r_2 r_3}{n_1 n_2 n_3}} c \sigma_{\min}(\mathcal{X}_\star).\end{aligned}$$

Putting these together, we have the advertised bound. \square

E.3 Sparse outliers

The following two lemmas are useful to control the sparse corruption term, of which the second lemma follows directly from translating [CLY21, Lemma 5] to the tensor case.

Lemma 11 ([YPCC16, Lemma 1] [CLY21, Lemma 6]). *Suppose that $\mathcal{S} \in \mathbb{R}^{m \times n}$ is α -sparse. Then one has*

$$\|\mathcal{S}\| \leq \alpha \sqrt{mn} \|\mathcal{S}\|_\infty, \quad \|\mathcal{S}\|_{2,\infty} \leq \sqrt{\alpha n} \|\mathcal{S}\|_\infty, \quad \text{and} \quad \|\mathcal{S}\|_{1,\infty} \leq \alpha n \|\mathcal{S}\|_\infty.$$

Lemma 12 ([CLY21, Lemma 5]). *Suppose that $\mathcal{Y} = \mathcal{X}_\star + \mathcal{S}_\star$ for some α -sparse \mathcal{S}_\star . Fix a tensor \mathcal{X} , and let $\mathcal{S} = \mathcal{T}_\zeta(\mathcal{Y} - \mathcal{X})$ where the threshold satisfies $\zeta \geq \|\mathcal{X} - \mathcal{X}_\star\|_\infty$. We then have*

$$\|\mathcal{S} - \mathcal{S}_\star\|_\infty \leq \|\mathcal{X} - \mathcal{X}_\star\|_\infty + \zeta \leq 2\zeta \tag{76}$$

and

$$\text{supp}(\mathcal{S}) \subseteq \text{supp}(\mathcal{S}_\star). \tag{77}$$

The relation (77) also implies that $\mathcal{S} - \mathcal{S}_\star$ is α -sparse.