

Nonconvex Low-Rank Tensor Completion from Noisy Data

Changxiao Cai* Gen Li† H. Vincent Poor* Yuxin Chen*

June 4, 2021

Abstract

We study a noisy tensor completion problem of broad practical interest, namely, the reconstruction of a low-rank tensor from highly incomplete and randomly corrupted observations of its entries. While a variety of prior work has been dedicated to this problem, prior algorithms either are computationally too expensive for large-scale applications, or come with sub-optimal statistical guarantees. Focusing on “incoherent” and well-conditioned tensors of a constant CP rank, we propose a two-stage nonconvex algorithm — (vanilla) gradient descent following a rough initialization — that achieves the best of both worlds. Specifically, the proposed nonconvex algorithm faithfully completes the tensor and retrieves all individual tensor factors within nearly linear time, while at the same time enjoying near-optimal statistical guarantees (i.e. minimal sample complexity and optimal estimation accuracy). The estimation errors are evenly spread out across all entries, thus achieving optimal ℓ_∞ statistical accuracy. We have also discussed how to extend our approach to accommodate asymmetric tensors. The insight conveyed through our analysis of nonconvex optimization might have implications for other tensor estimation problems.

Keywords: tensor completion, nonconvex optimization, gradient descent, spectral methods, entrywise statistical guarantees, minimaxity

Contents

1	Introduction and motivation	3
1.1	Tensor completion from noisy entries	3
1.2	Computational and statistical challenges	4
2	Algorithm and main results	5
2.1	A two-stage nonconvex algorithm	5
2.2	Main results	6
2.3	Numerical experiments	11
2.4	Notation	12
3	Initialization	13
3.1	Step 1: subspace estimation via a spectral method	13
3.2	Step 2: retrieval of low-rank tensor factors from the subspace estimate	13
3.2.1	Procedure	13
3.2.2	Intuition	14
3.3	Other alternatives?	14

Corresponding author: Yuxin Chen. A short version of this work has appeared in NeurIPS 2019 [CLPC19].

*Department of Electrical Engineering, Princeton University, Princeton, NJ 08544, USA; Email: {ccai, poor, yuxin.chen}@princeton.edu.

†Department of Electronic Engineering, Tsinghua University, Beijing, 10084, China; Email: g-li16@mails.tsinghua.edu.cn.

4	Related work	16
5	Analysis	17
5.1	Analysis for local convergence of GD	17
5.1.1	Preliminaries: gradient and Hessian calculation	17
5.1.2	Local strong convexity and smoothness	18
5.1.3	Leave-one-out gradient descent sequences	18
5.1.4	Key lemmas	19
5.2	Analysis for initialization: Part 1 (subspace estimation)	21
5.2.1	Key results	21
5.2.2	Leave-one-out sequences for subspace estimation	21
5.3	Analysis for initialization: Part 2 (retrieval of individual tensor factors)	22
5.3.1	Main results and leave-one-out sequences	22
5.3.2	Analysis	24
6	Discussion	27
A	Proofs for local convergence of GD	27
A.1	Proof of Lemma 5.1	28
A.1.1	Bounding α_4	28
A.1.2	Bounding α_1	29
A.1.3	Bounding α_2	31
A.1.4	Bounding α_3	31
A.1.5	Putting all this together	33
A.2	Proof of Lemma 5.2	33
A.3	Proof of Lemma 5.3	33
A.4	Proof of Lemma 5.4	36
A.5	Proof of Lemma 5.5	41
A.6	Proof of Lemma 5.6	45
B	Proofs for retrieving tensor components	45
B.1	Proof of Lemma 5.12	45
B.2	Proof of Lemma 5.13	47
B.3	Proof of Lemma B.1	49
B.3.1	Controlling F^τ	49
B.3.2	Controlling C^τ	51
B.4	Proof of Lemma 5.14	51
B.4.1	Controlling α_3	52
B.4.2	Controlling α_2	53
B.4.3	Controlling α_1	54
B.4.4	Combining α_1 , α_2 and α_3	55
B.5	Proof of Lemma B.2	55
B.5.1	Proximity of $M^{\tau,(m)}$ and $\widehat{M}^{\tau,(m)}$	55
B.5.2	Proximity of M^τ and $\widehat{M}^{\tau,(m)}$	56
B.6	Proof of Lemma 5.15	57
B.7	Proof of Lemma B.3	60
B.8	Proof of Lemma 5.16	61
B.8.1	Controlling β_1	62
B.8.2	Controlling β_2	63
B.8.3	Controlling β_3	63
B.8.4	Combining β_1 , β_2 and β_3	63
B.9	Proof of Lemma 5.17	64
B.10	Proof of Corollary 5.11	65
C	Proof of Corollary 2.9	65

D	Auxiliary lemmas	67
D.1	Statements of auxiliary lemmas	67
D.2	Proof of Lemma D.1	69
D.3	Proof of Lemma D.2 and Corollary D.3	70
D.3.1	Proof of Lemma D.2	70
D.3.2	Proof of Corollary D.3	72
D.4	Proof of Lemma D.4	73
D.5	Proof of Lemma D.5	75
D.6	Proof of Lemma D.6	76
D.7	Proof of Lemma D.7	77
D.8	Proof of Lemma D.8	77
D.9	Proof of Lemma D.9	77
D.10	Proof of Lemma D.10	78
E	Extension to asymmetric tensors	79
E.1	Problem settings	79
E.2	Algorithms	80
E.3	Numerical experiments	80
E.4	Analysis ideas	82
E.5	Proof of Lemma E.2	85

1 Introduction and motivation

1.1 Tensor completion from noisy entries

Estimation of low-complexity models from highly incomplete observations is a fundamental task that spans a diverse array of science and engineering applications. Arguably one of the most extensively studied problems of this kind is matrix completion, where one wishes to recover a low-rank matrix given only partial entries [DR16, CC18]. Moving beyond matrix-type data, a natural higher-order generalization is *low-rank tensor completion*, which aims to reconstruct a low-rank tensor when the vast majority of its entries are unseen. There is certainly no shortage of applications that motivate the investigation of tensor completion (e.g. personalized medicine [SN16, Paw19], medical imaging [GRY11, SHKM14, CZA+17], seismic data analysis [KSS13, EAHK13], multi-dimensional harmonic retrieval [CC14, YLW+17]). One concrete example in operations research arises when learning the preference of individual customers for a collection of products on the basis of historical transactions [FL19, MP20]. Given the limited availability of transaction data (e.g. each customer might only have purchased very few products before), it is crucial to exploit multi-way customer-product interactions (e.g. users’ browsing and searching histories) in order to better predict the likelihood of a customer purchasing a new product. Clearly, the presence of missing data and the need of exploiting multi-way structure result in the task of tensor completion. Additionally, tensor completion finds important applications in visual data in-painting [LMWY13, LYX17], where one wishes to reconstruct video data (or a sequence of images) from incomplete measurements. The video data consist of at least two spatial variables and one temporal variable, whose intrinsic connections are often modeled via certain low-complexity tensors.

For the sake of clarity, we phrase the problem formally before we proceed, focusing on a simple model that already captures the intrinsic difficulty of tensor completion in many aspects.¹ Imagine we are asked to estimate a symmetric order-three tensor² $\mathbf{T}^* \in \mathbb{R}^{d \times d \times d}$ from a small number of noisy entries

$$T_{j,k,l} = T_{j,k,l}^* + E_{j,k,l}, \quad \forall (j, k, l) \in \Omega, \quad (1)$$

where $T_{j,k,l}$ is the observed noisy entry at location (j, k, l) , $E_{j,k,l}$ stands for the associated noise, and $\Omega \subseteq \{1, \dots, d\}^3$ is a symmetric index subset to sample from. For notational simplicity, we set $\mathbf{T} = [T_{j,k,l}]_{1 \leq j,k,l \leq d}$

¹We focus on symmetric order-3 tensors primarily for simplicity of presentation. Many of our findings naturally extend to the more general case with asymmetric tensors of possibly higher order. Detailed discussions are deferred to Appendix E due to the space limits.

²Here, a tensor $\mathbf{T} \in \mathbb{R}^{d \times d \times d}$ is said to be symmetric if $T_{j,k,l} = T_{k,j,l} = T_{k,l,j} = T_{l,k,j} = T_{j,l,k} = T_{l,j,k}$ for all $1 \leq j, k, l \leq d$.

and $\mathbf{E} = [E_{j,k,l}]_{1 \leq j,k,l \leq d}$, with $T_{j,k,l} = E_{j,k,l} = 0$ for any $(j,k,l) \notin \Omega$. We adopt a *random sampling* model such that each index (j,k,l) ($j \leq k \leq l$) is included in Ω independently with probability p . In addition, we know *a priori* that the unknown tensor $\mathbf{T}^* \in \mathbb{R}^{d \times d \times d}$ is a superposition of r rank-one tensors (often termed canonical polyadic (CP) decomposition if r is minimal)

$$\mathbf{T}^* = \sum_{i=1}^r \mathbf{u}_i^* \otimes \mathbf{u}_i^* \otimes \mathbf{u}_i^*, \quad \text{or more concisely,} \quad \mathbf{T}^* = \sum_{i=1}^r \mathbf{u}_i^{*\otimes 3}, \quad (2)$$

where each $\mathbf{u}_i^* \in \mathbb{R}^d$ represents one of the r low-rank tensor components / factors. Here and throughout, for any vectors $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathbb{R}^d$, the tensor $\mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c}$ is a $d \times d \times d$ array whose (j,k,l) -th entry is given by $a_j b_k c_l$. The primary question is this: can we hope to faithfully estimate \mathbf{T}^* , as well as the individual tensor factors $\{\mathbf{u}_i^*\}_{1 \leq i \leq r}$, from the partially revealed entries (1), assuming that r is reasonably small?

1.2 Computational and statistical challenges

Even though tensor completion conceptually resembles matrix completion in various ways, it is considerably more challenging than the matrix counterpart. This is perhaps not surprising, given that a plethora of natural tensor problems (e.g. computing the spectral norm, finding the best low-rank approximation) are all notoriously hard [HL13]. As a notable example, while matrix completion is often efficiently solvable under nearly minimal sample complexity [CR09, Gro11], all polynomial-time algorithms developed so far for tensor completion — even in the noise-free case — require a sample size at least exceeding the order of $rd^{3/2}$, which is substantially larger than the degrees of freedom (i.e. rd) underlying the model (2). In fact, it is widely conjectured that there exists a large computational barrier away from the information-theoretic sampling limits [BM16].

With this fundamental gap in mind, the current paper focuses on the regime (in terms of the sample size) that enables reliable tensor completion in polynomial time. A variety of algorithms have been proposed that enjoy some sort of theoretical guarantees in (at least part of) this regime, including but not limited to spectral methods [MS18, CLC⁺20], sum-of-squares hierarchy [BM16, PS17], nonconvex algorithms [JO14, XY17], and also convex relaxation (based on proper unfolding) [GRY11, HMGW15, RPP13, GQ14]. While these are all polynomial-time algorithms, most of the computational complexities supported by prior theory remain prohibitively high when dealing with large-scale tensor data — a point that we shall elaborate on later. The only exception is the unfolding-based spectral method, which, however, fails to achieve exact recovery as the noise vanishes. This leads to a critical question:

Q1: *Is there any linear-time algorithm that is guaranteed to work for low-rank tensor completion?*

Going beyond such computational concerns, one might naturally wonder whether it is also possible for a fast algorithm to achieve a nearly un-improvable statistical accuracy in the presence of noise. Towards this end, intriguing stability guarantees have been established for sum-of-squares hierarchy in the noisy settings [BM16], although this paradigm is computationally expensive for large-scale data. The recent work [XYZ17] came up with a two-stage algorithm (i.e. a spectral method followed by tensor power iterations) for noisy tensor completion. Its estimation accuracy, however, falls short of achieving exact recovery in the absence of noise. This gives rise to another question of fundamental importance:

Q2: *Can we achieve near-optimal statistical accuracy without compromising computational efficiency?*

In this paper, we aim to address the above two questions by developing a nonconvex algorithm that achieves optimal computational efficiency and statistical accuracy all at once.

2 Algorithm and main results

2.1 A two-stage nonconvex algorithm

To address the above-mentioned challenges, a first impulse is to resort to the following least squares problem:

$$\underset{\mathbf{u}_1, \dots, \mathbf{u}_r \in \mathbb{R}^d}{\text{minimize}} \quad \sum_{j,k,l \in \Omega} \left(\left[\sum_{i=1}^r \mathbf{u}_i^{\otimes 3} \right]_{j,k,l} - T_{j,k,l} \right)^2, \quad (3)$$

or more concisely (up to proper re-scaling),

$$\underset{\mathbf{U} \in \mathbb{R}^{d \times r}}{\text{minimize}} \quad f(\mathbf{U}) := \frac{1}{6p} \left\| \mathcal{P}_\Omega \left(\sum_{i=1}^r \mathbf{u}_i^{\otimes 3} - \mathbf{T} \right) \right\|_F^2 \quad (4)$$

if we take $\mathbf{U} := [\mathbf{u}_1, \dots, \mathbf{u}_r] \in \mathbb{R}^{d \times r}$. Here, we denote by $\mathcal{P}_\Omega(\mathbf{T})$ the orthogonal projection of any tensor \mathbf{T} onto the subspace of tensors which vanish outside of the index set Ω . This optimization problem, however, is highly nonconvex (which involves minimizing a degree-6 polynomial), thus resulting in computational intractability in general.

Fortunately, not all nonconvex problems are as daunting to solve as they may seem. For example, recent years have seen a flurry of activity in low-rank matrix factorization via nonconvex optimization, which provably achieves optimal statistical accuracy and computational efficiency at once; see [CLC19] for an overview of recent advances. Motivated by this strand of work, we propose to solve (4) via a two-stage nonconvex paradigm, presented below in reverse order. The whole procedure is summarized in Algorithms 1-3.

Gradient descent (GD). Arguably one of the simplest optimization algorithms is gradient descent, which adopts a gradient update rule

$$\mathbf{U}^{t+1} = \mathbf{U}^t - \eta_t \nabla f(\mathbf{U}^t), \quad t = 0, 1, \dots \quad (5)$$

where η_t is the learning rate or the stepsize, and $\mathbf{U}^t \in \mathbb{R}^{d \times r}$ is the estimate in the t -th iteration. The main computational burden in each iteration lies in gradient evaluation, which, in this case, can be performed in time proportional to that taken to read the data.

Despite the simplicity of this algorithm, two critical issues stand out and might significantly affect its efficiency, which we shall bear in mind throughout the algorithmic and theoretical development.

(i) *Local stationary points and initialization.* As is well known, GD is guaranteed to find an approximate local stationary point, provided that the learning rates do not exceed the inverse Lipschitz constant of the gradient [Bub15]. There exist, however, local stationary points (e.g. saddle points or spurious local minima) that might fall short of the desired statistical properties. This requires us to properly avoid such undesired points, while retaining computational efficiency. To address this issue, one strategy is to first identify a rough initial guess within a local region surrounding the global solution (which often helps rule out bad local minima), in order to guarantee proper convergence of subsequent optimization procedures [LT17, JO14]. As a side remark, while careful initialization might not be crucial for several matrix recovery cases [CCFM19, GBW18, TV19], it does seem to be critical in various tensor problems [RM14]. We shall elucidate this point in Section 3.3.

(ii) *Learning rates and regularization.* Learning rates play a pivotal role in determining the convergence properties of GD. The challenge, however, is that the loss function (4) is overall not sufficiently smooth (i.e. its gradient often has an exceedingly large Lipschitz constant), and hence generic optimization theory recommends a pessimistically slow update rule (i.e. an extremely small learning rate) so as to guard against over-shooting. This, however, slows down the algorithm significantly, thus destroying the main computational advantage of GD (i.e. low per-iteration cost). With this issue in mind, prior literature suggests carefully designed regularization steps (e.g. proper projection, regularized loss functions) in order to improve the geometry of the optimization landscape [XY17]. In contrast, we argue that one is allowed to take a constant learning rate — which is as aggressive as it can possibly be — even without enforcing any regularization procedures.

Algorithm 1 Gradient descent for nonconvex tensor completion

- 1: Generate an initial estimate $\mathbf{U}^0 \in \mathbb{R}^{d \times r}$ via Algorithm 2.
 - 2: **for** $t = 0, 1, \dots, t_0 - 1$ **do**
 - 3: $\mathbf{U}^{t+1} = \mathbf{U}^t - \eta_t \nabla f(\mathbf{U}^t) = \mathbf{U}^t - \frac{\eta_t}{p} \mathcal{P}_\Omega(\sum_{i=1}^r (\mathbf{u}_i^t)^{\otimes 3} - \mathbf{T}) \times_1^{\text{seq}} \mathbf{U}^t \times_2^{\text{seq}} \mathbf{U}^t$, where \times_1^{seq} and \times_2^{seq} are defined in Section 2.4.
-

Algorithm 2 Spectral initialization for nonconvex tensor completion

- 1: Let $\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$ be the rank- r eigen-decomposition of

$$\mathbf{B} := \mathcal{P}_{\text{off-diag}}(\mathbf{A}\mathbf{A}^\top), \quad (6)$$

where $\mathbf{A} = \text{unfold}(p^{-1}\mathbf{T})$ is the mode-1 matricization of $p^{-1}\mathbf{T}$, and $\mathcal{P}_{\text{off-diag}}(\mathbf{Z})$ extracts out the off-diagonal entries of \mathbf{Z} .

- 2: **Output:** an initial estimate $\mathbf{U}^0 \in \mathbb{R}^{d \times r}$ on the basis of $\mathbf{U} \in \mathbb{R}^{d \times r}$ using Algorithm 3.
-

Initialization. Motivated by the above-mentioned issue (i), we develop a procedure that guarantees a reasonable initial estimate. In a nutshell, the proposed procedure consists of two steps:

- (a) Estimate the subspace spanned by the r low-rank tensor factors $\{\mathbf{u}_i^*\}_{1 \leq i \leq r}$ via a spectral method;
- (b) Disentangle individual low-rank tensor factors from this subspace estimate.

As we shall see momentarily, the total computational complexity of the proposed initialization is $O(pd^3)$ when $r = O(1)$, $\kappa = O(1)$ and $p \geq 1/d^2$ (where κ is a sort of “condition number” defined later), which is a linear-time algorithm. Note, however, that these two steps in the initialization procedure are relatively more complicated to describe. To improve the flow of the current paper, we postpone the details to Section 3. The readers can catch a glimpse of these procedures in Algorithms 2-3.

2.2 Main results

Encouragingly, the proposed nonconvex algorithm provably achieves the best of both worlds — in terms of statistical accuracy and computational efficiency — for a class of low-rank, well-conditioned, and “incoherent” problem instances. This subsection summarizes our main findings.

Before continuing, we note that one cannot hope to recover an arbitrary tensor from highly sub-sampled and arbitrarily corrupted entries. In order to enable provably valid recovery, the present paper focuses on a tractable model by imposing the following assumptions.

Definition 2.1 (Incoherence and well-conditionedness). *Define the incoherence parameters and the condition number of \mathbf{T}^* as follows*

$$\mu_0 := \frac{d^3 \|\mathbf{T}^*\|_\infty^2}{\|\mathbf{T}^*\|_F^2}, \quad (8a)$$

$$\mu_1 := \frac{d \|\mathbf{u}_i^*\|_\infty^2}{\|\mathbf{u}_i^*\|_2^2}, \quad (8b)$$

$$\mu_2 := \frac{d \langle \mathbf{u}_i^*, \mathbf{u}_j^* \rangle^2}{\|\mathbf{u}_i^*\|_2^2 \|\mathbf{u}_j^*\|_2^2}, \quad (8c)$$

$$\kappa := \frac{\max_i \|\mathbf{u}_i^*\|_2}{\min_i \|\mathbf{u}_i^*\|_2}. \quad (8d)$$

Remark 2.2. *Here, μ_0 , μ_1 and μ_2 are termed the incoherence parameters. Definitions (8a)-(8c) can be viewed as some sort of incoherence conditions for the tensor. For instance, when μ_0, μ_1 and μ_2 are small, these conditions say that (1) the energy of tensor \mathbf{T}^* is (nearly) evenly spread across all entries; (2) each factor \mathbf{u}_i^* is de-localized; (3) the factors $\{\mathbf{u}_i^*\}$ are nearly orthogonal to each other. Definition (8d) is concerned*

Algorithm 3 Retrieval of low-rank tensor factors from a given subspace estimate.

- 1: **Input:** number of restarts L , pruning threshold ϵ_{th} , subspace estimate $\mathbf{U} \in \mathbb{R}^{d \times r}$ given by Algorithm 2.
 - 2: **for** $\tau = 1, \dots, L$ **do**
 - 3: Generate an independent Gaussian vector $\mathbf{g}^\tau \sim \mathcal{N}(0, \mathbf{I}_d)$.
 - 4: $(\boldsymbol{\nu}^\tau, \lambda_\tau, \text{spec-gap}_\tau) \leftarrow \text{RETRIEVE-ONE-TENSOR-FACTOR}(\mathbf{T}, p, \mathbf{U}, \mathbf{g}^\tau)$.
 - 5: Generate tensor factor estimates $\{(\mathbf{w}^1, \lambda_1), \dots, (\mathbf{w}^r, \lambda_r)\} \leftarrow \text{PRUNE}(\{(\boldsymbol{\nu}^\tau, \lambda_\tau, \text{spec-gap}_\tau)\}_{\tau=1}^L, \epsilon_{\text{th}})$.
 - 6: **Output:** initial estimate $\mathbf{U}^0 = [\lambda_1^{1/3} \mathbf{w}^1, \dots, \lambda_r^{1/3} \mathbf{w}^r]$.
-

- 1: **function** RETRIEVE-ONE-TENSOR-FACTOR($\mathbf{T}, p, \mathbf{U}, \mathbf{g}$)
- 2: Compute

$$\boldsymbol{\theta} = \mathbf{U}\mathbf{U}^\top \mathbf{g} =: \mathcal{P}_{\mathbf{U}}(\mathbf{g}), \quad (7a)$$

$$\mathbf{M} = p^{-1} \mathbf{T} \times_3 \boldsymbol{\theta}, \quad (7b)$$

where \times_3 is defined in Section 2.4.

- 3: Let $\boldsymbol{\nu}$ be the leading singular vector of \mathbf{M} obeying $\langle \mathbf{T}, \boldsymbol{\nu}^{\otimes 3} \rangle \geq 0$, and set $\lambda = \langle p^{-1} \mathbf{T}, \boldsymbol{\nu}^{\otimes 3} \rangle$.
 - 4: **return** $(\boldsymbol{\nu}, \lambda, \sigma_1(\mathbf{M}) - \sigma_2(\mathbf{M}))$.
-

with the “well-conditionedness” of the tensor, meaning that each rank-1 component is of roughly the same size. In particular, we note that an assumption on pairwise correlation (i.e. a constraint on μ_2) is often assumed in the literature of tensor decomposition / factorization (e.g. [AGJ14, SLLC17, HZC20]).

For notational simplicity, we shall set

$$\mu := \max\{\mu_0, \mu_1, \mu_2\}. \quad (9)$$

Note that our theory allows μ to grow with the problem dimension d (in fact, μ can be as large as $d/\text{poly} \log(d)$).

Assumption 2.3 (Random noise). *Suppose that \mathbf{E} is a symmetric random tensor, where $\{E_{j,k,l}\}_{1 \leq j \leq k \leq l \leq d}$ (cf. (1)) are independently generated sub-Gaussian random variables with mean zero and variance $\text{Var}(E_{j,k,l}) \leq \sigma^2$.*

In addition, recognizing that there is a global permutational ambiguity issue (namely, one cannot distinguish $\mathbf{u}_1^*, \dots, \mathbf{u}_r^*$ from an arbitrary permutation of them), we introduce the following loss metrics to account for this ambiguity:

$$\text{dist}_{\text{F}}(\mathbf{U}, \mathbf{U}^*) := \min_{\mathbf{\Pi} \in \text{perm}_r} \|\mathbf{U}\mathbf{\Pi} - \mathbf{U}^*\|_{\text{F}}, \quad (10a)$$

$$\text{dist}_{\infty}(\mathbf{U}, \mathbf{U}^*) := \min_{\mathbf{\Pi} \in \text{perm}_r} \|\mathbf{U}\mathbf{\Pi} - \mathbf{U}^*\|_{\infty}, \quad (10b)$$

$$\text{dist}_{2,\infty}(\mathbf{U}, \mathbf{U}^*) := \min_{\mathbf{\Pi} \in \text{perm}_r} \|\mathbf{U}\mathbf{\Pi} - \mathbf{U}^*\|_{2,\infty}, \quad (10c)$$

where perm_r stands for the set of $r \times r$ permutation matrices. For notational simplicity, we also take

$$\lambda_{\min}^* := \min_{1 \leq i \leq r} \|\mathbf{u}_i^*\|_2^3 \quad \text{and} \quad \lambda_{\max}^* := \max_{1 \leq i \leq r} \|\mathbf{u}_i^*\|_2^3. \quad (11)$$

-
- 1: **function** PRUNE($\{(\boldsymbol{\nu}^\tau, \lambda_\tau, \text{spec-gap}_\tau)\}_{\tau=1}^L, \epsilon_{\text{th}}$)
 - 2: Set $\Theta = \{(\boldsymbol{\nu}^\tau, \lambda_\tau, \text{spec-gap}_\tau)\}_{\tau=1}^L$.
 - 3: **for** $i = 1, \dots, r$ **do**
 - 4: Choose $(\boldsymbol{\nu}^\tau, \lambda_\tau, \text{spec-gap}_\tau)$ from Θ with the largest spec-gap_τ ; set $\mathbf{w}^i = \boldsymbol{\nu}^\tau$ and $\lambda_i = \lambda_\tau$.
 - 5: Update $\Theta \leftarrow \Theta \setminus \{(\boldsymbol{\nu}^\tau, \lambda_\tau, \text{spec-gap}_\tau) \in \Theta : |\langle \boldsymbol{\nu}^\tau, \mathbf{w}^i \rangle| > 1 - \epsilon_{\text{th}}\}$.
 - 6: **return** $\{(\mathbf{w}^1, \lambda_1), \dots, (\mathbf{w}^r, \lambda_r)\}$.
-

With these notations in place, we are ready to present our main results. For simplicity of presentation, we shall start with the setting where $r, \mu, \kappa \asymp 1$.

Theorem 2.4. *Fix an arbitrary small constant $\delta > 0$. Suppose that $r, \kappa, \mu = O(1)$,*

$$p \geq c_0 \frac{\log^4 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^*} \leq c_1 \frac{\sqrt{p}}{d^{3/4} \log^2 d},$$

$$L = c_2 \quad \text{and} \quad \epsilon_{\text{th}} = c_3 \left(\frac{\log d}{d\sqrt{p}} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log^2 d}{p}} + \sqrt{\frac{\log d}{d}} \right)$$

for some sufficiently large constants $c_0, c_2 > 0$ and some sufficiently small constants $c_1, c_3 > 0$. The learning rate $\eta_t \equiv \eta$ is taken to be a constant obeying $0 < \eta \leq \lambda_{\min}^{*4/3} / (32\lambda_{\max}^{*8/3})$. Then with probability at least $1 - \delta$,

$$\text{dist}_{\text{F}}(\mathbf{U}^t, \mathbf{U}^*) \leq \left(C_1 \rho^t + C_2 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{\text{F}}, \quad (12a)$$

$$\text{dist}_{\infty}(\mathbf{U}^t, \mathbf{U}^*) \leq \text{dist}_{2,\infty}(\mathbf{U}^t, \mathbf{U}^*) \leq \left(C_3 \rho^t + C_4 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty} \quad (12b)$$

hold simultaneously for all $0 \leq t \leq t_0 = d^5$. Here, $0 < C_1, C_3, \rho < 1$ and $C_2, C_4 > 0$ are some absolute constants.

Remark 2.5. *The theorem holds unchanged if d^5 is replaced by d^c for an arbitrarily large constant $c > 0$.*

Remark 2.6. *The upper bound t_0 on the iteration count arises from the leave-one-out analysis when handling noisy observations. In short, the leave-one-out argument can only provide high-probability bounds for each iteration, thus requiring an upper bound on the iteration count if we desire a uniform bound across iterations. Note that in the noiseless case, our results and analysis hold for an arbitrarily large number of iterations.*

As an immediate consequence of Theorem 2.4, we obtain appealing ℓ_{∞} statistical guarantees for estimating tensor entries, which are previously rarely available (see Table 1). Specifically, let our tensor estimate in the t -th iteration be

$$\mathbf{T}^t := \sum_{i=1}^r \mathbf{u}_i^t \otimes \mathbf{u}_i^t \otimes \mathbf{u}_i^t, \quad \text{where } \mathbf{U}^t = [\mathbf{u}_1^t, \dots, \mathbf{u}_r^t] \in \mathbb{R}^{d \times r}. \quad (13)$$

Then our result is this:

Corollary 2.7. *Fix an arbitrarily small constant $\delta > 0$. Instate the assumptions of Theorem 2.4. Then with probability at least $1 - \delta$,*

$$\|\mathbf{T}^t - \mathbf{T}^*\|_{\text{F}} \lesssim \left(C_1 \rho^t + C_2 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{T}^*\|_{\text{F}}, \quad (14a)$$

$$\|\mathbf{T}^t - \mathbf{T}^*\|_{\infty} \lesssim \left(C_3 \rho^t + C_4 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{T}^*\|_{\infty} \quad (14b)$$

hold simultaneously for all $0 \leq t \leq t_0 = d^5$. Here, $0 < C_1, C_3, \rho < 1$ and $C_2, C_4 > 0$ are some absolute constants.

Several important implications are provided as follows. The discussion below assumes $\lambda_{\max}^* \asymp \lambda_{\min}^* \asymp 1$ for notational simplicity.

1. *Linear convergence.* In the absence of noise, the proposed algorithm converges linearly, namely, it provably attains ε accuracy within $O(\log(1/\varepsilon))$ iterations. Given the inexpensiveness of each gradient iteration, this algorithm can be viewed as a linear-time algorithm, which can almost be implemented as long as we can read the data. In the noisy setting, the algorithm reaches an appealing statistical accuracy within a logarithmic number of iterations.

	algorithm	sample complexity	computational complexity	ℓ_2 error (noisy)	ℓ_∞ error (noisy)	recovery type (noiseless)
our theory	spectral method + (vanilla) GD	$d^{1.5}$	pd^3	$\sigma\sqrt{\frac{d}{p}}$	$\sigma\sqrt{\frac{1}{p}}$	exact
[XYZ17]	spectral initialization + tensor power method	$d^{1.5}$	pd^3	$(\ \mathbf{T}^*\ _\infty + \sigma)\sqrt{\frac{d}{p}}$	n/a	approximate
[XY17]	spectral method + GD on manifold	$d^{1.5}$	poly(d)	n/a	n/a	exact
[MS18]	spectral method	$d^{1.5}$	d^3	n/a	n/a	approximate
[BM16]	sum-of-squares	$d^{1.5}$	d^{15}	$\frac{\ \mathbf{T}^*\ _F}{\sqrt{pd^{1.5}}} + \sigma d^{1.5}$	n/a	approximate
[PS17]	sum-of-squares	$d^{1.5}$	d^{10}	n/a	n/a	exact
[YZ16]	tensor nuclear norm	d	NP-hard	n/a	n/a	exact
[YZ17]	minimization	d	NP-hard	n/a	n/a	exact

Table 1: Comparison with prior theory for existing methods when $r, \mu, \kappa \asymp 1$ (neglecting logarithmic factors).

2. *Near-optimal sample complexity.* The fast convergence is guaranteed as soon as the sample size exceeds the order of $d^{3/2}\text{poly}\log d$. This matches the minimal sample complexity — modulo some logarithmic factor — known so far for any polynomial-time algorithm.
3. *Near-optimal statistical accuracy.* The proposed algorithm converges geometrically fast to a point with Euclidean error $O(\sigma\sqrt{(d\log d)/p})$. This matches the lower bound established in [XYZ17, Theorem 5] up to some logarithmic factor, thus justifying the statistical optimality of the proposed nonconvex algorithm.
4. *Entrywise estimation accuracy.* In addition to the Euclidean statistical guarantees, we have also established an entrywise error bound, which, to the best of our knowledge, has not been established in any of the prior work. When t is sufficiently large, the iterates reach an entrywise error bound $O(\sigma\sqrt{(\log d)/p})$. This entrywise error bound is about an order of \sqrt{d} times smaller than the above ℓ_2 error bound, thereby implying that the estimation errors are evenly spread out across all entries.
5. *Noise size.* The above theory operates in the regime where $\sigma \lesssim \sqrt{\frac{p}{d^{3/2}}}$ (modulo some log factor). Given that we have $\|\mathbf{T}^*\|_\infty \asymp d^{-3/2}$ in this case, our noise size constraint can be equivalently written as (up to some log factor)

$$\frac{\sigma}{\|\mathbf{T}^*\|_\infty} \lesssim \sqrt{pd^{3/2}}. \quad (15)$$

Since the sampling rate needs to satisfy $p \gg d^{-3/2}$, this condition essentially allows the typical size of each noise component to be considerably larger than the size of the corresponding entry of the truth, which covers a broad range of practical scenarios.

6. *Implicit regularization.* One appealing feature of our finding is the simplicity of the iterative refinement stage of the algorithm. All of the above statistical and computational benefits hold for vanilla gradient descent (when properly initialized). This should be contrasted with prior work (e.g. [XY17]) that relies on extra regularization terms to stabilize the optimization landscape. In principle, vanilla gradient descent implicitly constrains itself within a region of well-conditioned landscape, thus enabling fast convergence without explicit regularization.
7. *No need of sample splitting.* The theory developed herein does not require fresh samples in each iteration. We note that sample splitting has been frequently adopted in other context primarily to simplify mathematical analysis. Nevertheless, it typically does not exploit the data in an efficient manner (i.e. each data sample is used only once), thus resulting in the need of a much larger sample size in practice.

We shall take a moment to discuss the merits of our approach in comparison to prior work. One of the best-known polynomial-time algorithms is the degree-6 level of the sum-of-squares (SoS) hierarchy, which seems to

match the computationally feasible limit in terms of the sample complexity [BM16]. However, this approach has a well-documented limitation in that it involves solving a semidefinite program of dimensions $d^3 \times d^3$, which requires enormous storage and computation power. The work [MS18] alleviates this computational burden by resorting to a clever unfolding-based spectral algorithm; it is a nearly linear-time procedure that enables near-minimal sample complexity (among polynomial-time algorithms), although it does not achieve exact recovery even in the absence of noise. The two-stage algorithm developed by [XYZ17] — which is based on spectral initialization followed by tensor power methods — shares similar advantages and drawbacks as [MS18]. Further, the recent work [XY17] proposes a polynomial-time nonconvex algorithm based on gradient descent over Grassmann manifold (with a properly regularized objective function), which is an extension of the nonconvex matrix completion algorithm proposed by [KMO10a, KMO10b] to tensor data. The theory provided in [XY17], however, does not provide explicit computational complexities. The recent work [SY19] attempts tensor estimation via an interesting algorithm adapted from collaborative filtering and investigates both ℓ_2 and ℓ_∞ estimation accuracy. This approach, however, does not guarantee exact recovery in the absence of noise. We summarize and compare several prior results in Table 1 (omitting logarithmic factors).

Thus far, we have concentrated on the low-rank, well-conditioned, and incoherent case. Our main theory can be extended to cover a broader class of scenarios, as stated below.

Theorem 2.8. *Fix an arbitrary small constant $\delta > 0$. Suppose that $\kappa \asymp 1$,*

$$p \geq c_0 \frac{\mu^4 r^4 \log^4 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^*} \leq c_1 \frac{\sqrt{p}}{\mu r^{3/2} d^{3/4} \log^2 d}, \quad r \leq c_2 \left(\frac{d}{\mu^6 \log^6 d} \right)^{1/6},$$

$$L = c_3 r^{2\kappa^2} \log^{3/2} r \quad \text{and} \quad \epsilon_{\text{th}} = c_4 \left(\frac{\mu r \log d}{d\sqrt{p}} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{rd \log^2 d}{p}} + \sqrt{\frac{\mu r \log d}{d}} \right)$$

for some sufficiently large constants $c_0, c_3 > 0$ and some sufficiently small constants $c_1, c_2, c_4 > 0$. The learning rate $\eta_t \equiv \eta$ is taken to be a constant obeying $0 < \eta \leq \lambda_{\min}^{*4/3} / (32\lambda_{\max}^{*8/3})$. Then with probability at least $1 - \delta$,

$$\text{dist}_{\text{F}}(\mathbf{U}^t, \mathbf{U}^*) \leq \left(C_1 \rho^t + C_2 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{\text{F}} \quad (16a)$$

$$\text{dist}_{\infty}(\mathbf{U}^t, \mathbf{U}^*) \leq \text{dist}_{2,\infty}(\mathbf{U}^t, \mathbf{U}^*) \leq \left(C_3 \rho^t + C_4 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty} \quad (16b)$$

hold simultaneously for all $0 \leq t \leq t_0 = d^5$. Here, $0 < C_1, C_3, \rho < 1$ and $C_2, C_4 > 0$ are some absolute constants.

Corollary 2.9. *Fix an arbitrarily small constant $\delta > 0$. Instate the assumptions of Theorem 2.8. Then with probability at least $1 - \delta$,*

$$\|\mathbf{T}^t - \mathbf{T}^*\|_{\text{F}} \lesssim \left(C_1 \rho^t + C_2 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{T}^*\|_{\text{F}}, \quad (17a)$$

$$\|\mathbf{T}^t - \mathbf{T}^*\|_{\infty} \lesssim \left(C_3 \rho^t + C_4 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{\mu^3 r d \log d}{p}} \right) \|\mathbf{T}^*\|_{\infty} \quad (17b)$$

hold simultaneously for all $0 \leq t \leq t_0 = d^5$. Here, $0 < C_1, C_3, \rho < 1$ and $C_2, C_4 > 0$ are some absolute constants.

Remark 2.10. *Clearly, Theorem 2.8 and Corollary 2.9 subsume Theorem 2.4 and Corollary 2.7 as a special case respectively.*

Remark 2.11. *Our theorems require the rank r to not exceed $o(d^{1/6})$, which, we believe, is an artifact of the current nonconvex analysis (particularly for the initialization stage). For instance, our local convergence*

analysis is built upon strong convexity and smoothness, which holds only within a sufficiently small neighborhood surrounding the truth; given that the diameter of this neighborhood is no more than $o(1/r)$, our analysis requires an initial guess with higher accuracy than expected, thus leading to our rank constraint. It might be possible to improve the rank dependency via more refined analysis, and we leave it to future investigation.

In a nutshell, this theorem reveals intriguing theoretical support (including both ℓ_F and $\ell_{2,\infty}$ bounds) for more general settings. Assuming that the condition number $\kappa \asymp 1$, the nonconvex algorithm we propose is guaranteed to succeed in polynomial time. Note, however, that our theoretical dependency (including both sample and computational complexities) on the rank r and the incoherence parameter μ are likely loose and sub-optimal. In addition, if κ is allowed to grow with d , then the current theory requires a large number of restart attempts during the initialization stage, resulting in a very high computational burden. Improving these aspects, however, calls for a much more refined analysis framework, which we leave for future investigation.

2.3 Numerical experiments

We carry out a series of numerical experiments to corroborate our theoretical findings. Before proceeding, recall that Theorem 2.8 only guarantees successful recovery with probability $1 - \delta$ for some small constant δ ; this means that we shall not anticipate a very high success rate (e.g. $1 - O(d^{-5})$) as in the matrix recovery case. As we shall make clear shortly, this happens mainly because the initialization stage works only with probability $1 - \delta$, where the uncertainty largely depends on the random vectors $\{\mathbf{g}^\tau\}_{1 \leq \tau \leq L}$. With this observation in mind, we recommend the following modification to improve the empirical success rate:

- Run Algorithm 2 independently for $t_{\text{init}} = 5$ times to obtain multiple initial estimates (denoted by $\mathbf{U}_{[1]}^0, \dots, \mathbf{U}_{[t_{\text{init}}]}^0$); select the one achieving the smallest empirical loss, namely

$$\mathbf{U}_{\text{best}}^0 = \arg \min_{\mathbf{U} \in \{\mathbf{U}_{[i]}^0\}_{1 \leq i \leq t_{\text{init}}}} f(\mathbf{U}). \quad (18)$$

- Run Algorithm 1 with the initial point \mathbf{U}^0 set to be $\mathbf{U}_{\text{best}}^0$.

The final estimates for the low-rank factor and the whole tensor are denoted respectively by

$$\widehat{\mathbf{U}} = \mathbf{U}^{t_0} \quad \text{and} \quad \widehat{\mathbf{T}} = \sum_{i=1}^r \mathbf{u}_i^{t_0} \otimes \mathbf{u}_i^{t_0} \otimes \mathbf{u}_i^{t_0}, \quad (19)$$

where $\mathbf{U}^{t_0} = [\mathbf{u}_1^{t_0}, \dots, \mathbf{u}_r^{t_0}] \in \mathbb{R}^{d \times r}$ is the iterate returned by Algorithm 1, with t_0 the total number of gradient iterations. In the sequel, we generate the true tensor $\mathbf{T}^* = \sum_{1 \leq i \leq r} \mathbf{u}_i^* \otimes^3$ randomly in such a way that $\mathbf{u}_i^* \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. The learning rates are taken to be $\eta_t \equiv 0.2$ unless otherwise noted.

We start with numerical convergence rates of our algorithm in the absence of noise. Set $d = 100$, $r = 4$, $p = 0.1$, $L = 16$ and $\epsilon_{\text{th}} = 0.4$. Figure 1(a) the numerical estimation errors vs. iteration count t in a typical Monte Carlo trial. Here, four kinds of estimation errors are reported: (1) the relative Frobenius norm error $\frac{\text{dist}_F(\mathbf{U}^t, \mathbf{U}^*)}{\|\mathbf{U}^*\|_F}$; (2) the relative $\|\cdot\|_{2,\infty}$ error $\frac{\text{dist}_{2,\infty}(\mathbf{U}^t, \mathbf{U}^*)}{\|\mathbf{U}^*\|_{2,\infty}}$; (3) the relative Frobenius norm error $\frac{\|\mathbf{T}^t - \mathbf{T}^*\|_F}{\|\mathbf{T}^*\|_F}$; (4) the relative ℓ_∞ error $\frac{\|\mathbf{T}^t - \mathbf{T}^*\|_\infty}{\|\mathbf{T}^*\|_\infty}$. Here, $\mathbf{T}^t = \sum_{i=1}^r \mathbf{u}_i^t \otimes \mathbf{u}_i^t \otimes \mathbf{u}_i^t$ with $\mathbf{U}^t = [\mathbf{u}_1^t, \dots, \mathbf{u}_r^t]$. For all these metrics, the numerical estimation errors decay geometrically fast.

Next, we study the phase transition (in terms of the success rates for exact recovery) in the noise-free settings. Set $d = 100$, $r = 4$, $L = 16$ and $\epsilon_{\text{th}} = 0.4$. For the sake of comparisons, we also report the numerical performance of the tensor power method (TPM) followed by gradient descent. When running the tensor power method, we set both the number of iterations and the restart number to be 16. Each trial is claimed to succeed if the relative ℓ_2 error obeys $\frac{\text{dist}_F(\widehat{\mathbf{U}}, \mathbf{U}^*)}{\|\mathbf{U}^*\|_F} \leq 0.01$. Figure 1(b) plots the empirical success rates over 100 independent Monte Carlo trials. As can be seen, our initialization algorithm outperforms the tensor power method.

The third series of experiments is concerned with the dependence of the success rate on the rank r . Let us set $p = rd^{-3/2} \log^2 d$, $L = r^2$ and $\epsilon_{\text{th}} = 0.4$, and the success recovery criterion is the same as above. Figure 1(c) depicts the empirical success rates (over 100 independent Monte Carlo trials) as the rank r varies. As

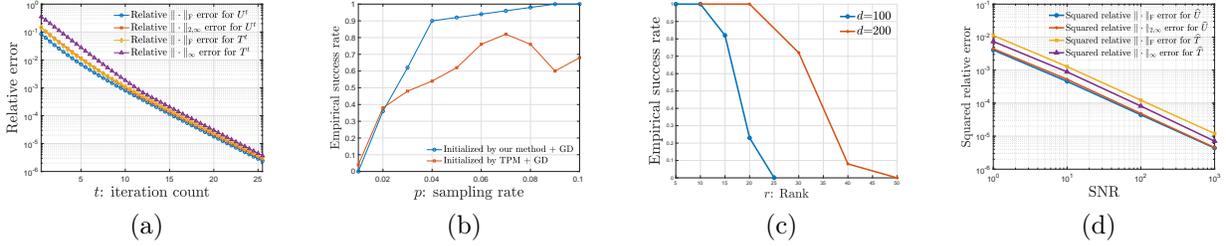


Figure 1: (a) relative errors of the estimates \mathbf{U}^t and \mathbf{T}^t vs. iteration count t for noiseless tensor completion, where $d = 100$, $r = 4$, $p = 0.1$; (b) empirical success rate vs. sampling rate, where $d = 100$, $r = 4$; (c) empirical success rate vs. rank, where $p = rd^{-3/2} \log^2 d$; (d) squared relative errors vs. SNR for noisy settings, where $d = 100$, $r = 4$ and $p = 0.1$. Each point in (b), (c) and (d) is averaged over 100 independent Monte Carlo trials.

can be seen from the plots, the proposed algorithm is able to achieve exact reconstruction as long as the rank r is sufficiently small compared to d . The plausible range of r , however, seems and seems to be larger than our theoretic requirement $r = o(d^{1/6})$. This, once again, suggests the need of future investigation to pin down the best possible dependency on r .

Finally, we consider the numerical estimation accuracy of our algorithm. Take $t_0 = 100$, $d = 100$, $r = 4$, $p = 0.1$, $L = 16$ and $\epsilon_{\text{th}} = 0.4$. Define the signal-to-noise ratio (SNR) to be $\text{SNR} = \frac{\|\mathbf{T}^*\|_{\text{F}}^2/d^3}{\sigma^2}$. We report in Figure 1(d) three types of squared relative errors (namely, $\frac{\text{dist}_{\text{F}}^2(\hat{\mathbf{U}}, \mathbf{U}^*)}{\|\mathbf{U}^*\|_{\text{F}}^2}$, $\frac{\text{dist}_{2,\infty}^2(\hat{\mathbf{U}}, \mathbf{U}^*)}{\|\mathbf{U}^*\|_{2,\infty}^2}$ and $\frac{\|\hat{\mathbf{T}} - \mathbf{T}^*\|_{\infty}^2}{\|\mathbf{T}^*\|_{\infty}^2}$) vs. SNR. Figure 1(d) illustrates that all three types of relative squared errors scale inversely proportional to the SNR (since the slope in the figure is roughly -1), which is consistent with our statistical guarantees.

2.4 Notation

Before proceeding, we gather a few notations that will be used throughout this paper. First of all, for any matrix $\mathbf{M} \in \mathbb{R}^{d \times d}$, we let $\|\mathbf{M}\|$ and $\|\mathbf{M}\|_{\text{F}}$ denote the operator norm (or the spectral norm) and the Frobenius norm of \mathbf{M} , respectively, and let $\mathbf{M}_{i,:}$ and $\mathbf{M}_{:,i}$ denote the i -th row and i -th column, respectively. In addition, $\lambda_1(\mathbf{M}) \geq \lambda_2(\mathbf{M}) \geq \dots \geq \lambda_d(\mathbf{M})$ denote the eigenvalues of \mathbf{M} and $\sigma_1(\mathbf{M}) \geq \sigma_2(\mathbf{M}) \geq \dots \geq \sigma_d(\mathbf{M})$ denote the singular values of \mathbf{M} .

For any tensor $\mathbf{T} \in \mathbb{R}^{d \times d \times d}$, let $\mathbf{T}_{i,:,:} \in \mathbb{R}^{d \times d}$ denote the mode-1 i -slice with entries $(\mathbf{T}_{i,:,:})_{j,k} = T_{i,j,k}$, and $\mathbf{T}_{:,i,:}$ and $\mathbf{T}_{:,:,i}$ are defined in a similar way. For any tensors $\mathbf{T}, \mathbf{R} \in \mathbb{R}^{d \times d \times d}$, the inner product is defined as $\langle \mathbf{T}, \mathbf{R} \rangle := \sum_{j,k,l} T_{j,k,l} R_{j,k,l}$. The Frobenius norm of \mathbf{T} is defined as $\|\mathbf{T}\|_{\text{F}} := \sqrt{\langle \mathbf{T}, \mathbf{T} \rangle}$. For any vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$, we define the vector products of a tensor $\mathbf{T} \in \mathbb{R}^{d \times d \times d}$ — denoted by $\mathbf{T} \times_3 \mathbf{u} \in \mathbb{R}^{d \times d}$ and $\mathbf{T} \times_1 \mathbf{u} \times_2 \mathbf{v} \in \mathbb{R}^d$ — such that

$$[\mathbf{T} \times_3 \mathbf{u}]_{ij} := \sum_{1 \leq k \leq d} T_{i,j,k} u_k, \quad 1 \leq i, j \leq d; \quad (20a)$$

$$[\mathbf{T} \times_1 \mathbf{u} \times_2 \mathbf{v}]_k := \sum_{1 \leq i, j \leq d} T_{i,j,k} u_i v_j, \quad 1 \leq k \leq d. \quad (20b)$$

The products $\mathbf{T} \times_2 \mathbf{u} \in \mathbb{R}^{d \times d}$, $\mathbf{T} \times_3 \mathbf{u} \in \mathbb{R}^{d \times d}$, $\mathbf{T} \times_1 \mathbf{u} \times_3 \mathbf{v} \in \mathbb{R}^d$ and $\mathbf{T} \times_2 \mathbf{u} \times_3 \mathbf{v} \in \mathbb{R}^d$ are defined in a similar manner. For any $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_r] \in \mathbb{R}^{d \times r}$ and $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_r] \in \mathbb{R}^{d \times r}$, we further define

$$\mathbf{T} \times_1^{\text{seq}} \mathbf{U} \times_2^{\text{seq}} \mathbf{V} := [\mathbf{T} \times_1 \mathbf{u}_i \times_2 \mathbf{v}_i]_{1 \leq i \leq r} \in \mathbb{R}^{d \times r}. \quad (21)$$

In addition, the operator norm of \mathbf{T} is defined as

$$\|\mathbf{T}\| := \sup_{\mathbf{u}, \mathbf{v}, \mathbf{w} \in \mathbb{S}^{d-1}} \langle \mathbf{T}, \mathbf{u} \otimes \mathbf{v} \otimes \mathbf{w} \rangle, \quad (22)$$

where $\mathbb{S}^{d-1} := \{\mathbf{u} \in \mathbb{R}^d \mid \|\mathbf{u}\|_2 = 1\}$ indicates the unit sphere in \mathbb{R}^d .

Further, $f(n) \lesssim g(n)$ or $f(n) = O(g(n))$ means that $|f(n)/g(n)| \leq C_1$ for some constant $C_1 > 0$; $f(n) \gtrsim g(n)$ means that $|f(n)/g(n)| \geq C_2$ for some constant $C_2 > 0$; $f(n) \asymp g(n)$ means that $C_1 \leq |f(n)/g(n)| \leq C_2$ for some constants $C_1, C_2 > 0$; $f(n) = o(g(n))$ means that $\lim_{n \rightarrow \infty} f(n)/g(n) = 0$. In addition, $f(n) \ll g(n)$ means that $f(n) \leq c_1 g(n)$ for some sufficiently small constant $c_1 > 0$, and $f(n) \gg g(n)$ means that $f(n) \geq c_2 g(n)$ for some sufficiently large constant $c_2 > 0$.

3 Initialization

This section presents formal details of the proposed two-step initialization, accompanied by some intuition. Recall that the proposed initialization procedure consists of two steps, which we discuss separately.

3.1 Step 1: subspace estimation via a spectral method

The spectral algorithm is often applied in conjunction with simple “unfolding” (or “matricization”) to estimate the *subspace* spanned by the r factors $\{\mathbf{u}_i^*\}_{1 \leq i \leq r}$. This strategy is partly motivated by prior approaches developed for covariance estimation with missing data [Lou14, MS18, CLC⁺20]. We provide a brief introduction below.

Let

$$\mathbf{A} = \text{unfold}^{1 \times 2} \left(\frac{1}{p} \mathbf{T} \right) \in \mathbb{R}^{d \times d^2}, \quad \text{or more concisely } \mathbf{A} = \text{unfold} \left(\frac{1}{p} \mathbf{T} \right) \in \mathbb{R}^{d \times d^2} \quad (23)$$

be the mode-1 matricization of $p^{-1} \mathbf{T}$ (namely, $\frac{1}{p} T_{i,j,k} = A_{i,(j-1)d+k}$ for any $1 \leq i, j, k \leq d$) [KB09]. The rationale of this step is that: under our model, the unfolded matrix \mathbf{A} obeys

$$\mathbb{E}[\mathbf{A}] = \text{unfold}(\mathbf{T}^*) = \sum_{i=1}^r \mathbf{u}_i^* (\mathbf{u}_i^* \otimes \mathbf{u}_i^*)^\top =: \mathbf{A}^*, \quad (24)$$

whose column space is precisely the span of $\{\mathbf{u}_i^*\}_{1 \leq i \leq r}$. This motivates one to estimate the r -dimensional column space of $\mathbb{E}[\mathbf{A}]$ from \mathbf{A} . Towards this, a natural strategy is to look at the principal subspace of $\mathbf{A} \mathbf{A}^\top$. However, the diagonal entries of $\mathbf{A} \mathbf{A}^\top$ bear too much influence on the principal directions and need to be properly down-weighted. The current paper chooses to work with the principal subspace of the following matrix that zeros out all diagonal components:

$$\mathbf{B} := \mathcal{P}_{\text{off-diag}}(\mathbf{A} \mathbf{A}^\top), \quad (25)$$

where $\mathcal{P}_{\text{off-diag}}(\mathbf{Z})$ extracts out the off-diagonal entries of a squared matrix \mathbf{Z} . If we let $\mathbf{U} \in \mathbb{R}^{d \times r}$ be an orthonormal matrix whose columns are the top- r eigenvectors of \mathbf{B} , then \mathbf{U} serves as our subspace estimate. See Algorithm 2 for a summary of the procedure.

3.2 Step 2: retrieval of low-rank tensor factors from the subspace estimate

3.2.1 Procedure

As it turns out, it is possible to obtain rough (but reasonable) estimates of all individual low-rank tensor factors $\{\mathbf{u}_i^*\}_{1 \leq i \leq r}$ — up to global permutation — given a reliable subspace estimate \mathbf{U} . This is in stark contrast to the low-rank matrix recovery case, where there exists some global rotational ambiguity that prevents us from disentangling the r factors of interest.

We begin by describing how to retrieve *one* tensor factor from the subspace estimate — a procedure summarized in RETRIEVE-ONE-TENSOR-FACTOR(). Let us generate a random vector from the provided subspace \mathbf{U} (which has orthonormal columns), that is,

$$\boldsymbol{\theta} = \underbrace{\mathbf{U} \mathbf{U}^\top}_{\text{projection of } \mathbf{g} \text{ onto } \mathbf{U}} \mathbf{g}, \quad \mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d). \quad (26)$$

The rescaled tensor data $p^{-1}\mathbf{T}$ is then transformed into a matrix via proper “projection” along this random direction $\boldsymbol{\theta}$, namely,

$$\mathbf{M} = \frac{1}{p}\mathbf{T} \times_3 \boldsymbol{\theta} \in \mathbb{R}^{d \times d}. \quad (27)$$

Our estimate for a tensor factor is then given by $\lambda^{1/3}\boldsymbol{\nu}$, where $\boldsymbol{\nu}$ is the leading singular vector of \mathbf{M} obeying $\langle \mathbf{T}, \boldsymbol{\nu}^{\otimes 3} \rangle \geq 0$, and λ is taken as $\lambda = \langle p^{-1}\mathbf{T}, \boldsymbol{\nu}^{\otimes 3} \rangle$. Informally, $\boldsymbol{\nu}$ reflects the direction of the component \mathbf{u}_i^* that exhibits the largest correlation with the random direction $\boldsymbol{\theta}$, and λ forms an estimate of the corresponding size $\|\mathbf{u}_i^*\|_2$.

A challenge remains, however, as there are oftentimes more than one tensor factors to estimate. To address this issue, we propose to re-run the aforementioned procedure multiple times, so as to ensure that we get to retrieve each tensor factor of interest at least once. We will then apply a careful pruning procedure (i.e. PRUNE()) to remove redundancy.

3.2.2 Intuition

To develop some intuition about the above procedure, consider the “heuristic” case where $\boldsymbol{\theta} = \mathbf{U}^*(\mathbf{U}^{*\top}\mathbf{U}^*)^{-1}\mathbf{U}^{*\top}\mathbf{g}$, namely, the idealistic scenario where the subspace estimate \mathbf{U} is accurate. Averaging out the randomness in the sampling pattern and the noise, we see that the expected projected matrix (27) takes the following form:

$$\mathbb{E}[\mathbf{M} \mid \boldsymbol{\theta}] = \mathbf{T}^* \times_3 \boldsymbol{\theta} = \sum_{i=1}^r \langle \boldsymbol{\theta}, \mathbf{u}_i^* \rangle \mathbf{u}_i^* \mathbf{u}_i^{*\top}.$$

As a result, in the incoherent case where $\{\mathbf{u}_j^*\}$ are nearly orthogonal to each other, the leading singular vector of $\mathbb{E}[\mathbf{M} \mid \boldsymbol{\theta}]$ — and hence that of \mathbf{M} (i.e. \mathbf{w}) — is expected to be reasonably close to the factor \mathbf{u}_i^* that enjoys the largest projected coefficient. In other words, we expect

$$\boldsymbol{\nu} \approx \frac{1}{\|\mathbf{u}_i^*\|_2} \mathbf{u}_i^*, \quad \text{where } i = \arg \max_{1 \leq j \leq r} |\langle \boldsymbol{\theta}, \mathbf{u}_j^* \rangle|. \quad (28)$$

In the mean time, armed with (28) and the incoherence assumption (such that \mathbf{u}_i^* and \mathbf{u}_j^* are nearly orthogonal for $i \neq j$), one might have

$$\lambda = \langle \mathbf{T}^*, \boldsymbol{\nu}^{\otimes 3} \rangle \approx \frac{1}{\|\mathbf{u}_i^*\|_2^3} \langle \mathbf{T}^*, \mathbf{u}_i^{*\otimes 3} \rangle \approx \frac{1}{\|\mathbf{u}_i^*\|_2^3} \langle \mathbf{u}_i^{*\otimes 3}, \mathbf{u}_i^{*\otimes 3} \rangle = \|\mathbf{u}_i^*\|_2^3, \quad (29)$$

thus explaining our choice of λ in the proposed procedure. These arguments hint at the ability of our procedure in retrieving one tensor factor in each round.

The above intuitive argument, however, does not explain why we need to first project a random vector \mathbf{g} onto the (approximate) column space of \mathbf{U}^* . While we won’t go into detailed calculations here, we remark in passing a crucial high variability issue: without proper projection, the perturbation incurred by both the missing data and the noise might far exceed the strength of the true signal. As a result, it is advised to first project the data onto the desired subspace, in the hope of amplifying the signal-to-noise ratio.

3.3 Other alternatives?

The careful reader may naturally wonder whether a careful initialization is pivotal in achieving fast convergence. While a thorough answer to this has yet to be developed, we shall point out some alternatives that seem sub-optimal in both theory and practice. To simplify the presentation, the current subsection focuses on the rank-1 noiseless case, where

$$\mathbf{T}^* = \mathbf{u}^{*\otimes 3}, \quad \mathbf{T} = \frac{1}{p}\mathcal{P}_\Omega(\mathbf{T}^*), \quad \|\mathbf{u}^*\|_2 = 1. \quad (30)$$

Since the decision variable is now a d -dimensional vector, we shall employ the conventional notation \mathbf{u}^t to represent \mathbf{U}^t .

Random initialization. We find it instrumental to begin with the population-level analysis, which corresponds to the scenario with no missing data and noise ($p = 1$ and $\sigma = 0$). A little calculation gives

$$\mathbb{E} [\mathbf{u}^1 | \mathbf{u}^0] = \mathbb{E} [\mathbf{u}^0 - \eta \nabla f(\mathbf{u}^0) | \mathbf{u}^0] = (1 - \eta \|\mathbf{u}^0\|_2^4) \mathbf{u}^0 + \eta \langle \mathbf{u}^0, \mathbf{u}^* \rangle^2 \mathbf{u}^*. \quad (31)$$

As an immediate consequence, the expected correlation between the next iterate and the truth obeys

$$\mathbb{E} [\langle \mathbf{u}^1, \mathbf{u}^* \rangle | \mathbf{u}^0] = \{1 - \eta \|\mathbf{u}^0\|_2^4 + \eta \langle \mathbf{u}^0, \mathbf{u}^* \rangle \|\mathbf{u}^*\|_2^2\} \langle \mathbf{u}^0, \mathbf{u}^* \rangle.$$

This means that if \mathbf{u}^0 and \mathbf{u}^* are positively correlated and if the initial guess \mathbf{u}^0 is sufficiently small,³ then one has

$$\mathbb{E} [\langle \mathbf{u}^1, \mathbf{u}^* \rangle | \mathbf{u}^0] \approx (1 + \eta \langle \mathbf{u}^0, \mathbf{u}^* \rangle \|\mathbf{u}^*\|_2^2) \langle \mathbf{u}^0, \mathbf{u}^* \rangle; \quad (32)$$

a similar recursion holds for \mathbf{u}^t . As a result, the GD iterates are expected to get increasingly more aligned with the truth, at least at the population level. Caution needs to be exercised, however, that this population-level analysis alone fails to capture what is happening in the finite-sample case. In what follows, we point out potential issues with random initialization.

Consider the case where \mathbf{u}^0 is generated as a vector of i.i.d. Gaussian random variables. Suppose that \mathbf{u}^0 and \mathbf{u}^* are positively correlated and that $\|\mathbf{u}^0\|_2$ is sufficiently small. It is easily seen that, with high probability, the expected increment is on the order of (cf. (32))

$$\mathbb{E} [\langle \mathbf{u}^1, \mathbf{u}^* \rangle | \mathbf{u}^0] - \langle \mathbf{u}^0, \mathbf{u}^* \rangle \approx \eta \langle \mathbf{u}^0, \mathbf{u}^* \rangle^2 \|\mathbf{u}^*\|_2^2 \lesssim \frac{\eta \text{poly log}(d)}{d} \|\mathbf{u}^*\|_2^4 \|\mathbf{u}^0\|_2^2, \quad (33)$$

which could be quite small as it depends quadratically on the current correlation $\langle \mathbf{u}^0, \mathbf{u}^* \rangle$.

If we were to hope that the favorable population-level analysis captures more or less the finite-sample dynamics, we would need to ensure that the variability of the gradient update is well-controlled. Towards this, let us compute the variance of $\langle \mathbf{u}^1, \mathbf{u}^* \rangle$, assuming that $\frac{\|\mathbf{u}^0\|_\infty}{\|\mathbf{u}^0\|_2} \asymp \frac{\|\mathbf{u}^*\|_\infty}{\|\mathbf{u}^*\|_2} \asymp \frac{\text{poly log}(d)}{\sqrt{d}}$:

$$\begin{aligned} \text{Var} (\langle \mathbf{u}^1, \mathbf{u}^* \rangle | \mathbf{u}^0) &\asymp \text{Var} \left(\frac{\eta}{p} \sum_{1 \leq j, k, l \leq d} (\chi_{jkl} - p) (u_j^0 u_k^0 u_l^0 - u_j^* u_k^* u_l^*) u_j^0 u_k^0 u_l^* \right) \\ &\asymp \frac{\eta^2}{p} \sum_{1 \leq j, k, l \leq d} (u_j^0 u_k^0 u_l^0 + u_j^* u_k^* u_l^*)^2 (u_j^0 u_k^0 u_l^*)^2 \asymp \frac{\eta^2 \text{poly log}(d)}{pd^3} \|\mathbf{u}^*\|_2^8 \|\mathbf{u}^0\|_2^4. \end{aligned}$$

In other words, the typical size of the variability of $\langle \mathbf{u}^1, \mathbf{u}^* \rangle$ is about the order of $\frac{\eta \text{poly log}(d)}{\sqrt{pd^3}} \|\mathbf{u}^*\|_2^4 \|\mathbf{u}^0\|_2^2$, which dominates (in fact, is order-of-magnitudes larger than) the mean increment (33) unless

$$p \gtrsim \frac{\text{poly log}(d)}{d}. \quad (34)$$

The sample size corresponding to (34) is, however, considerably larger than the computation limit $p \asymp \frac{\text{poly log}(d)}{d^{1.5}}$. The presence of a large variance implies highly volatile dynamics of randomly initialized GD, thus casting doubt on its efficiency in the most challenging sample-starved regime.

In summary, the main issue stems from the quadratic dependence of the expected increment (33) on the correlation $\langle \mathbf{u}^0, \mathbf{u}^* \rangle$, which can be exceedingly small if \mathbf{u}^0 is randomly initialized.

Initialization via the tensor power method (TPM). Another alternative for initialization is the tensor power method, which has recently gained popularity in the context of learning latent-variable models [AGH⁺14, AGJ17]. Nevertheless, the TPM (with random initialization) suffers from the same high-volatility issue as randomly initialized GD. The argument for this would be nearly identical to the one presented

³In fact, if a random initialization \mathbf{u}^0 is not small, then one can easily show that, with high probability, the ℓ_2 norm of \mathbf{u}^t is going to drop geometrically fast at the beginning.

above, and is hence omitted. Instead, we invoke a perturbation analysis result in [AGH⁺14, Theorem 5.1] to illustrate the insufficiency of the TPM.

Recall that $\frac{1}{p}\mathbf{T} = \mathbf{T}^* + (\frac{1}{p}\mathbf{T} - \mathbf{T}^*)$. A critical issue is that the perturbation bound in [AGH⁺14, Theorem 5.1] requires the tensor perturbation to be exceedingly small, namely,

$$\|\frac{1}{p}\mathbf{T} - \mathbf{T}^*\| \lesssim 1/d. \quad (35)$$

This, however, cannot possibly hold if the sample size is merely $p \asymp \frac{\text{poly log}(d)}{d^{1.5}}$ (in which case one only expects a spectral norm bound on the order of $\|p^{-1}\mathbf{T} - \mathbf{T}^*\| \lesssim \frac{1}{\text{poly log}(d)}$ shown in Corollary D.3 even in the absence of noise). In light of all this, existing stability analysis of the TPM does not imply either sample efficiency or computational efficiency.

4 Related work

One of the most natural ideas for solving tensor completion is to first unfold the tensor data into matrices, followed by proper convex relaxation commonly adopted for low-rank matrix completion. Given that there are more than one ways to matricize a tensor, several prior work has explored the design of matrix norms that can exploit the tensor structure more effectively [THK10, GRY11, LMWY13, RPP13, LFC⁺16, MHWG14]. Such algorithms have been robustified to enable reliable recovery against sparse outliers as well [GQ14]. For the most part, however, such unfolding-based convex relaxation necessarily incur loss of structural information, which is particularly severe when handling odd-order tensors. The sample complexity developed for this paradigm is often sub-optimal vis-a-vis the computational limits (namely, minimal sample complexity achievable by polynomial-time algorithms).

Motivated by the above sub-optimality issue, [YZ16, YZ17] proposed to minimize instead the tensor nuclear norm subject to data constraints, which provably allows for reduced sample complexity. The issue, however, is that computing the tensor nuclear norm itself is already computationally intractable, thus limiting its applicability to even moderate-dimensional problems. Similar findings have also been discovered for tensor atomic norm minimization [DBBG19]. When restricted to polynomial-time algorithms, the best statistical guarantees are often attained via convex relaxation tailored to the sum-of-squares hierarchy [BM16]; the resulting computational cost, however, remains prohibitively high for practical large-scale problems. Another matrix nuclear norm minimization algorithm has been proposed based on promoting certain structures on certain factor matrices [LSC⁺14]. Developing statistical guarantees is, however, not the focal point of this work.

Moving beyond convex relaxation, a number of prior papers have developed nonconvex algorithms for tensor completion, examples including iterative hard thresholding [RSS17], alternating minimization [JO14, WAA16, XHYS15], tensor SVD [ZA17], optimization on manifold [XY17, KM16, Ste16], proximal average algorithm with nonconvex regularizer [Yao18], and block coordinate decent [JHZ⁺16, XY13]. When it comes to the model considered herein, these algorithms either lack optimal statistical guarantees, or come with a computational cost that is significantly higher than a linear-time algorithm.

The algorithm and theory that we develop are largely inspired by the recent advances of nonconvex optimization algorithms for low-rank matrix recovery problems [KMO10a, KMO10b, CLS15, CC17, SL16, YPCC16, CW15]. The main theoretical tool — the leave-one-out analysis — is a powerful technique that has proved successful in various other statistical problems [EK15, CFMW19, AFWZ17, MWCC17, ZB18, CCFM19, CCF⁺19, LZT19, CFMY19, DC18, PW19]. There are several major differences between the analysis of nonconvex tensor completion and that of nonconvex matrix recovery. For instance, our initialization scheme is substantially more complicated than the matrix recovery counterpart, thus requiring much more sophisticated analysis; in addition, the local convergence stage of tensor completion does not suffer from rotational ambiguity (which often appears in nonconvex matrix completion), and hence we only need to handle permutational ambiguity.

In addition, the current paper focuses on non-adaptive uniform random sampling. If there is freedom in designing the sampling mechanism, then one can often expect improved performance; see [KS13, Zha19] as examples. Fundamental criteria that enable perfect low-CP-rank tensor completion have been studied in [AW17].

Tensor completion is simply a special example of the tensor recovery literature. There is a large body of results tackling various other tensor recovery and estimation problems, including but not limited to tensor decomposition [Kol01, KB09, AGH⁺14, AGJ14, TS15, KOKC13, HSS16, GHJY15, ZKOM18, SDFL⁺17, SLLC17, GM17], tensor SVD and factorization [ZX18, KBHH13, ZA17], and tensor regression and sketching [RSS17, HZC20, CRY19, HWW⁺19]. The algorithmic ideas explored in this paper might have implications for these tensor-related problems as well.

5 Analysis

In this section, we outline the proof of Theorem 2.8. The proof of Corollary 2.9 is deferred to Appendix C. The analysis is divided into three parts:

- In Section 5.1, we show that given an initial estimate sufficiently close to the ground truth, vanilla gradient descent converges linearly. These are formalized in Lemmas 5.3 and 5.6.
- Sections 5.2-5.3 provide statistical guarantees for the two steps of the initialization procedure; see Theorems 5.9.
- Under the assumptions of Theorem 2.8, one can see that the initialization satisfies the requirement of linear convergence of vanilla gradient descent. Therefore, Theorem 2.8 immediately follows from the results in Sections 5.1-5.3.

5.1 Analysis for local convergence of GD

In this section, we demonstrate that: if the initialization is reasonably good, then vanilla gradient descent converges linearly to a solution with the desired statistical accuracy. We postpone the analysis for initialization to Sections 5.2-5.3 for convenience of presentation.

5.1.1 Preliminaries: gradient and Hessian calculation

First of all, using our notation \times^{seq} defined in (21), we can write

$$\nabla f(\mathbf{U}) = \frac{1}{p} \mathcal{P}_\Omega \left(\sum_{1 \leq i \leq r} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* - \mathbf{E} \right) \times_1^{\text{seq}} \mathbf{U} \times_2^{\text{seq}} \mathbf{U}. \quad (36)$$

Next, we find it convenient to define an auxiliary loss function $f_{\text{clean}}(\mathbf{U}) : \mathbb{R}^{d \times r} \rightarrow \mathbb{R}_+$ that corresponds to the noiseless case:

$$f_{\text{clean}}(\mathbf{U}) = \frac{1}{6p} \left\| \mathcal{P}_\Omega \left(\sum_{1 \leq i \leq r} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right) \right\|_{\text{F}}^2. \quad (37)$$

The gradient of f_{clean} w.r.t. \mathbf{u}_s ($1 \leq s \leq r$) is thus given by

$$\nabla_{\mathbf{u}_s} f_{\text{clean}}(\mathbf{U}) = \frac{1}{p} \mathcal{P}_\Omega \left(\sum_{1 \leq i \leq r} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right) \times_1 \mathbf{u}_s \times_2 \mathbf{u}_s, \quad 1 \leq s \leq r, \quad (38)$$

and hence one can write

$$\nabla f_{\text{clean}}(\mathbf{U}) = \frac{1}{p} \mathcal{P}_\Omega \left(\sum_{1 \leq i \leq r} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right) \times_1^{\text{seq}} \mathbf{U} \times_2^{\text{seq}} \mathbf{U}. \quad (39)$$

This clearly satisfies

$$\nabla f(\mathbf{U}) = \nabla f_{\text{clean}}(\mathbf{U}) - \frac{1}{p} \mathcal{P}_\Omega(\mathbf{E}) \times_1^{\text{seq}} \mathbf{U} \times_2^{\text{seq}} \mathbf{U}. \quad (40)$$

Moreover, direct algebraic manipulations give that: for any matrix $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_r] \in \mathbb{R}^{d \times r}$,

$$\begin{aligned} \text{vec}(\mathbf{V})^\top \nabla^2 f_{\text{clean}}(\mathbf{U}) \text{vec}(\mathbf{V}) &= \frac{1}{3p} \left\| \mathcal{P}_\Omega \left(\sum_{1 \leq s \leq r} \mathbf{u}_s \otimes \mathbf{u}_s \otimes \mathbf{v}_s + \mathbf{u}_s \otimes \mathbf{v}_s \otimes \mathbf{u}_s + \mathbf{v}_s \otimes \mathbf{u}_s \otimes \mathbf{u}_s \right) \right\|_{\text{F}}^2 \\ &\quad + \frac{2}{p} \left\langle \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{u}_s^{\otimes 3} - \mathbf{T}^* \right), \sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{v}_s \otimes \mathbf{u}_s \right\rangle, \end{aligned} \quad (41)$$

where $\text{vec}(\mathbf{V})$ denotes the vectorization of \mathbf{V} .

5.1.2 Local strong convexity and smoothness

At the heart of our analysis is a crucial geometric property of the objective function, that is, the noiseless loss function f_{clean} behaves like a locally strongly convex and smooth function. This fact, which is formally stated in the following lemma, is the key enabler of fast local convergence of vanilla GD.

Lemma 5.1 (Local strong convexity and smoothness). *Suppose that the sample complexity and the rank satisfy*

$$p \geq c_0 \max \left\{ \frac{\log^3 d}{d^{3/2}}, \frac{\mu^2 r^2 \log d}{d^2} \right\} \quad r \leq c_1 \sqrt{\frac{d}{\mu}} \quad (42)$$

for some sufficiently large (resp. small) constant $c_0 > 0$ (resp. $c_1 > 0$). Then with probability greater than $1 - O(d^{-10})$,

$$\frac{1}{2} \lambda_{\min}^{*4/3} \|\mathbf{V}\|_{\text{F}}^2 \leq \text{vec}(\mathbf{V})^\top \nabla^2 f_{\text{clean}}(\mathbf{U}) \text{vec}(\mathbf{V}) \leq 4 \lambda_{\max}^{*4/3} \|\mathbf{V}\|_{\text{F}}^2 \quad (43)$$

holds simultaneously for all $\mathbf{V} \in \mathbb{R}^{d \times r}$ and all $\mathbf{U} \in \mathbb{R}^{d \times r}$ obeying

$$\|\mathbf{U} - \mathbf{U}^*\|_{\text{F}} \leq \delta \|\mathbf{U}^*\|_{\text{F}} \quad \text{and} \quad \|\mathbf{U} - \mathbf{U}^*\|_{2,\infty} \leq \delta \|\mathbf{U}^*\|_{2,\infty}. \quad (44)$$

Here, $\delta \leq c_2/(\mu^{3/2}r)$ for some sufficiently small constant $c_2 > 0$.

Proof. See Appendix A.1. ■

In order to invoke Lemma 5.1, one needs to make sure that the decision matrix \mathbf{U} of interest (e.g. \mathbf{U}^t in the GD sequence) satisfies the condition (44). This, however, is a fairly stringent condition, as it requires \mathbf{U} to be close to the truth in every single row.

5.1.3 Leave-one-out gradient descent sequences

Motivated by the analytical framework developed for low-rank matrix recovery [MWCC17, CLL19], we introduce the following leave-one-out sequences, which play a crucial role in guaranteeing that the entire trajectory $\{\mathbf{U}^t\}_{t \geq 0}$ satisfies the condition (44) as required in Lemma 5.1.

Specifically, we define for each $1 \leq m \leq d$ the following auxiliary loss function:

$$f^{(m)}(\mathbf{U}) \triangleq \frac{1}{6p} \left\| \mathcal{P}_{\Omega_{-m}} \left(\sum_{1 \leq s \leq r} \mathbf{u}_s^{\otimes 3} - \mathbf{T}^* - \mathbf{E} \right) \right\|_{\text{F}}^2 + \frac{1}{6} \left\| \mathcal{P}_m \left(\sum_{1 \leq s \leq r} \mathbf{u}_s^{\otimes 3} - \mathbf{T}^* \right) \right\|_{\text{F}}^2, \quad (45)$$

where

- \mathcal{P}_{Ω_m} : the projection onto the subspace of tensors supported on $\{(i, j, k) \in \Omega : i = m \text{ or } j = m \text{ or } k = m\}$;
- $\mathcal{P}_{\Omega_{-m}}$: the projection onto the subspace of tensors supported on $\{(i, j, k) \in \Omega : i \neq m \text{ and } j \neq m \text{ and } k \neq m\}$;
- \mathcal{P}_m : the projection onto the subspace of tensors supported on $\{(i, j, k) \in [d]^3 : i = m \text{ or } j = m \text{ or } k = m\}$.

In words, this function is obtained by replacing all data at locations $\{(i, j, k) \in [d]^3 : i = m \text{ or } j = m \text{ or } k = m\}$ by their expected values, thus removing all randomness associated with this location subset. The gradient of $f^{(m)}(\mathbf{U})$ w.r.t. \mathbf{u}_s ($1 \leq s \leq r$) can be computed as:

$$\nabla_{\mathbf{u}_s} f^{(m)}(\mathbf{U}) = \frac{1}{p} \mathcal{P}_{\Omega_{-m}} \left(\sum_{1 \leq s \leq r} \mathbf{u}_s^{\otimes 3} - \mathbf{T}^* - \mathbf{E} \right) \times_1 \mathbf{u}_s \times_2 \mathbf{u}_s + \mathcal{P}_m \left(\sum_{1 \leq s \leq r} \mathbf{u}_s^{\otimes 3} - \mathbf{T}^* \right) \times_1 \mathbf{u}_s \times_2 \mathbf{u}_s. \quad (46)$$

We then denote by $\{\mathbf{U}^{t,(m)}\}_{t \geq 0}$ the iterative sequence obtained by running gradient descent w.r.t. the leave-one-out loss $f^{(m)}(\cdot)$; see Algorithm 4. By construction, as long as $\mathbf{U}^{0,(m)}$ is independent of the sampling locations and the noise associated with the locations $\{(i, j, k) \in \Omega : i = m \text{ or } j = m \text{ or } k = m\}$ (which holds true as detailed momentarily), then the entire trajectory $\{\mathbf{U}^{t,(m)}\}_{t \geq 0}$ becomes statistically independent of such randomness. This is a crucial property that allows us to decouple the complicated statistical dependency.

Algorithm 4 The m -th leave-one-out sequence

- 1: Generate an initial estimate $\mathbf{U}^{0,(m)}$ via Algorithm 5.
 - 2: **for** $t = 0, 1, \dots, t_0 - 1$ **do**
 - 3: $\mathbf{U}^{t+1,(m)} = \mathbf{U}^{t,(m)} - \eta_t \nabla f^{(m)}(\mathbf{U}^{t,(m)})$.
-

5.1.4 Key lemmas

The proof for local linear convergence of GD is inductive in nature, which proceeds on the basis of the following set of inductive hypotheses. As we shall see in Corollary 5.11 in Section 5.3, this set of inductive hypotheses — modulo some global permutation — is valid with high probability when $t = 0$. In order to simplify presentation, we remove the consideration of the global permutation factor throughout this section (namely, we assume that the following holds for $\mathbf{U}^0 \mathbf{\Pi}^0$ with some permutation matrix $\mathbf{\Pi}^0 \in \mathbb{R}^{r \times r}$ obeying $\mathbf{\Pi}^0 = \mathbf{I}$). Our inductive hypotheses are summarized as follows:

Key hypotheses for the gradient update stage:

$$\|\mathbf{U}^t - \mathbf{U}^*\|_{\text{F}} \leq \left(C_1 \rho^t \mathcal{E}_{\text{local}} + C_2 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{\text{F}}; \quad (47a)$$

$$\|\mathbf{U}^t - \mathbf{U}^*\|_{2,\infty} \leq \left(C_3 \rho^t \mathcal{E}_{\text{local}} + C_4 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty}; \quad (47b)$$

$$\|\mathbf{U}^t - \mathbf{U}^{t,(m)}\|_{\text{F}} \leq \left(C_5 \rho^t \mathcal{E}_{\text{local}} + C_6 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty}; \quad (47c)$$

$$\|(\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:}\|_2 \leq \left(C_7 \rho^t \mathcal{E}_{\text{local}} + C_8 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty}; \quad (47d)$$

for some quantity $\mathcal{E}_{\text{local}} > 0$ (depending possibly on μ and r) and some constants $C_1, \dots, C_8 > 0$. These exist a few straightforward consequences of the hypotheses (47), which we record in the following lemma.

Lemma 5.2. *Assume that the hypotheses (47) hold, then we have*

$$\|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\text{F}} \leq \left(2C_1 \rho^t \mathcal{E}_{\text{local}} + 2C_2 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{\text{F}}, \quad (48)$$

$$\|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{2,\infty} \leq \left((C_3 + C_5) \rho^t \mathcal{E}_{\text{local}} + (C_4 + C_6) \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty}. \quad (49)$$

Proof. See Appendix A.2. ■

Our proof for the hypotheses (47) is inductive in nature: we would like to show that if the hypotheses in (47) hold for the t -th iteration, then they continue to be valid for the $(t + 1)$ -th iteration. We shall justify each of the above hypotheses inductively through the following lemmas.

Lemma 5.3. *Suppose that*

$$p \geq c_0 \frac{\mu^3 r^2 \log^3 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^*} \leq c_1 \frac{\sqrt{p}}{\mu^{3/2} r \sqrt{d \log d}}, \quad \text{and} \quad r \leq c_2 \sqrt{\frac{d}{\mu}}$$

for some sufficiently large constant $c_0 > 0$ and some sufficiently small constant $c_1, c_2 > 0$. Assume that the hypotheses (47) hold for the t -th iteration and $\mathcal{E}_{\text{local}} \leq c_3 / (\mu^{3/2} r)$ for some sufficiently small constant $c_3 > 0$. Then with probability at least $1 - O(d^{-10})$,

$$\|\mathbf{U}^{t+1} - \mathbf{U}^*\|_{\text{F}} \leq \left(C_1 \rho^{t+1} \mathcal{E}_{\text{local}} + C_2 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{\text{F}}, \quad (50)$$

provided that $0 < \eta \leq \lambda_{\min}^{*4/3}/(32\lambda_{\max}^{*8/3})$, $1 - (\lambda_{\min}^{*4/3}/5)\eta \leq \rho < 1$, and C_2 is sufficiently large.

Proof. See Appendix A.3. ■

Lemma 5.4. *Suppose that*

$$p \geq c_0 \frac{\mu^3 r^2 \log^3 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^*} \leq c_1 \frac{\sqrt{p}}{\mu^{3/2} r \sqrt{d \log d}}, \quad \text{and} \quad r \leq c_2 \sqrt{\frac{d}{\mu}}$$

for some sufficiently large constant $c_0 > 0$ and some sufficiently small constant $c_1, c_2 > 0$. Assume that the hypotheses (47) hold for the t -th iteration and $\mathcal{E}_{\text{local}} \leq c_3/(\mu^{3/2}r)$ for some sufficiently small constant $c_3 > 0$. Then with probability at least $1 - O(d^{-10})$, one has

$$\|\mathbf{U}^{t+1,(m)} - \mathbf{U}^{t+1}\|_{\text{F}} \leq \left(C_5 \rho^{t+1} \mathcal{E}_{\text{local}} + C_6 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty}, \quad (51)$$

provided that $0 < \eta \leq \lambda_{\min}^{*4/3}/(32\lambda_{\max}^{*8/3})$, $1 - (\lambda_{\min}^{*4/3}/5)\eta \leq \rho < 1$ and C_6 is sufficiently large.

Proof. See Appendix A.4. ■

Lemma 5.5. *Suppose that*

$$p \geq c_0 \frac{\mu^3 r^2 \log^3 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^*} \leq c_1 \frac{\sqrt{p}}{\mu^{3/2} r \sqrt{d \log d}}, \quad \text{and} \quad r \leq c_2 \sqrt{\frac{d}{\mu}}$$

for some sufficiently large constant $c_0 > 0$ and some sufficiently small constant $c_1, c_2 > 0$. Assume that the hypotheses (47) hold for the t -th iteration and $\mathcal{E}_{\text{local}} \leq c_3/(\mu^{3/2}r)$ for some sufficiently small constant $c_3 > 0$. Then with probability at least $1 - O(d^{-10})$, one has

$$\|(\mathbf{U}^{t+1,(m)} - \mathbf{U}^*)_{m,:}\|_2 \leq \left(C_7 \rho^{t+1} \mathcal{E}_{\text{local}} + C_8 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty}, \quad (52)$$

provided that $0 < \eta \leq \lambda_{\min}^{*4/3}/(32\lambda_{\max}^{*8/3})$, $1 - (\lambda_{\min}^{*4/3}/5)\eta \leq \rho < 1$, C_7 and C_8 are sufficiently large.

Proof. See Appendix A.5. ■

Lemma 5.6. *Suppose that*

$$p \geq c_0 \frac{\mu^3 r^2 \log^3 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^*} \leq c_1 \frac{\sqrt{p}}{\mu^{3/2} r \sqrt{d \log d}}, \quad \text{and} \quad r \leq c_2 \sqrt{\frac{d}{\mu}}$$

for some sufficiently large constant $c_0 > 0$ and some sufficiently small constant $c_1, c_2 > 0$. Assume that the hypotheses (47) hold for the t -th iteration and $\mathcal{E}_{\text{local}} \leq c_3/(\mu^{3/2}r)$ for some sufficiently small constant $c_3 > 0$. Then with probability at least $1 - O(d^{-10})$, one has

$$\|\mathbf{U}^{t+1} - \mathbf{U}^*\|_{2,\infty} \leq \left(C_3 \rho^{t+1} \mathcal{E}_{\text{local}} + C_4 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty}, \quad (53)$$

provided that $0 < \eta \leq \lambda_{\min}^{*4/3}/(32\lambda_{\max}^{*8/3})$, $1 - (\lambda_{\min}^{*4/3}/5)\eta \leq \rho < 1$, $C_3/(C_5 + C_7)$ and $C_4/(C_6 + C_8)$ are both sufficiently large.

Proof. See Appendix A.6. ■

The proofs of the above key lemmas are postponed to Appendix A.

5.2 Analysis for initialization: Part 1 (subspace estimation)

5.2.1 Key results

The aim of this subsection is to demonstrate that the subspace estimate \mathbf{U} computed by Algorithm 2 is sufficiently close to the space spanned by the true tensor factors. Given that the columns of $\mathbf{U}^* = [\mathbf{u}_1^*, \dots, \mathbf{u}_r^*]$ are in general not orthogonal to each other, we shall define $\mathbf{U}_{\text{orth}}^* \in \mathbb{R}^{d \times r}$ as follows (obtained by proper orthonormalization) :

$$\mathbf{U}_{\text{orth}}^* := \mathbf{U}^* (\mathbf{U}^{*\top} \mathbf{U}^*)^{-\frac{1}{2}}. \quad (54)$$

This matrix $\mathbf{U}_{\text{orth}}^*$ reflects the rank- r principal subspace of $\mathbf{A}^* \mathbf{A}^{*\top} = \sum_i \|\mathbf{u}_i^*\|_2^4 \mathbf{u}_i^* \mathbf{u}_i^{*\top}$, where we recall that $\mathbf{A}^* \in \mathbb{R}^{d \times d^2}$ is the mode-1 matricization of \mathbf{T}^* . In addition, we define the rotation matrix

$$\mathbf{R} := \arg \min_{\mathbf{Q} \in \mathcal{O}^{r \times r}} \|\mathbf{U} \mathbf{Q} - \mathbf{U}_{\text{orth}}^*\|_{\text{F}}, \quad (55)$$

where $\mathcal{O}^{r \times r}$ stands for the set of $r \times r$ orthonormal matrices. This can be viewed as the global rotation matrix that best aligns the two subspaces represented by \mathbf{U} and $\mathbf{U}_{\text{orth}}^*$ respectively.

Equipped with the above notation, we can invoke [CLC⁺20, Corollary 1] to arrive at the following lemma, which upper bounds the distance between our subspace estimate \mathbf{U} and the ground truth $\mathbf{U}_{\text{orth}}^*$.

Lemma 5.7. *There exist some universal constants $c_0, c_1, c_2 > 0$ such that if*

$$p \geq c_0 \frac{\mu^2 r \log^2 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^*} \leq c_1 \frac{\sqrt{p}}{d^{3/4} \sqrt{\log d}}, \quad \text{and} \quad r \leq c_2 \sqrt{\frac{d}{\mu}},$$

then with probability $1 - O(d^{-10})$, the subspace estimate \mathbf{U} computed by Algorithm 2 obeys

$$\|\mathbf{U} \mathbf{R} - \mathbf{U}_{\text{orth}}^*\| \lesssim \mathcal{E}_{\text{se}}, \quad (56a)$$

$$\|\mathbf{U} \mathbf{R} - \mathbf{U}_{\text{orth}}^*\|_{2, \infty} \lesssim \mathcal{E}_{\text{se}} \sqrt{\frac{\mu r}{d}}, \quad (56b)$$

where $\mathbf{U}_{\text{orth}}^*$ and \mathbf{R} are defined respectively in (54) and (55), and

$$\mathcal{E}_{\text{se}} := \frac{\mu^2 r \log d}{d^{3/2} p} + \sqrt{\frac{\mu^2 r \log d}{d^2 p}} + \frac{\sigma^2}{\lambda_{\min}^{*2}} \frac{d^{3/2} \log d}{p} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} + \frac{\mu r}{d}. \quad (57)$$

In a nutshell, Lemma 5.7 asserts that: under our sample size, noise and rank conditions, Algorithm 2 produces reliable estimates of the subspace spanned by the low-rank tensor factors $\{\mathbf{u}_i^*\}_{1 \leq i \leq r}$. The theorem quantifies the subspace distance in terms of both the spectral norm and $\|\cdot\|_{2, \infty}$, where the latter bound often reflects a considerably stronger sense of proximity compared to the former one.

As it turns out, in order to facilitate analysis for the subsequent stages, we need to introduce certain leave-one-out sequences as well, which we detail in the next subsection.

5.2.2 Leave-one-out sequences for subspace estimation

The key idea of the leave-one-out analysis is to create auxiliary leave-one-out sequences that are (1) independent of a small fraction of the data; (2) sufficiently close to the true estimates. We introduce the following auxiliary tensor and $d \times d^2$ -dimensional matrix for each $1 \leq m \leq d$:

$$\mathbf{T}^{(m)} := \mathcal{P}_{\Omega_{-m}}(\mathbf{T}) + p \mathcal{P}_m(\mathbf{T}^*) \in \mathbb{R}^{d \times d \times d}, \quad (58)$$

$$\mathbf{A}^{(m)} := \text{mode-1 matricization of } \frac{1}{p} \mathbf{T}^{(m)}. \quad (59)$$

By construction, $\mathbf{T}^{(m)}$ and $\mathbf{A}^{(m)}$ are independent of $\mathcal{P}_{\Omega_m}(\mathbf{E})$, where we recall that

$$\Omega_{-m} := \{(i, j, k) \in \Omega : i \neq m \text{ and } j \neq m \text{ and } k \neq m\}, \quad (60)$$

$$\Omega_m := \{(i, j, k) \in \Omega : i = m \text{ or } j = m \text{ or } k = m\}. \quad (61)$$

We are now ready to introduce the auxiliary leave-one-out procedure for subspace estimation. Similar to the matrix \mathbf{B} in Algorithm 2 (whose eigenspace serves as an estimate of the column space of \mathbf{U}^*), we define an auxiliary matrix $\mathbf{B}^{(m)} \in \mathbb{R}^{d \times d}$ as follows:

$$\mathbf{B}^{(m)} = \mathcal{P}_{\text{off-diag}}(\mathbf{A}^{(m)} \mathbf{A}^{(m)\top}), \quad (62)$$

where $\mathcal{P}_{\text{off-diag}}(\cdot)$ (as already defined in Section 3.1) extracts out off-diagonal entries from a matrix. The rationale is simple: it can be easily verified that

$$\mathbb{E}[\mathbf{B}^{(m)}] = \mathbf{B}^* - \mathcal{P}_{\text{diag}}(\mathbf{B}^*), \quad \mathbf{B}^* := \mathbf{A}^* \mathbf{A}^{*\top}, \quad (63)$$

where $\mathcal{P}_{\text{diag}}(\cdot)$ extracts out the diagonal entries of the matrix. This gives hope that the eigenspace of $\mathbf{B}^{(m)}$ is also a reliable estimate of the column space of \mathbf{U}^* , provided that the diagonal entries of \mathbf{B}^* are sufficiently small. Consequently, we shall compute $\mathbf{U}^{0,(m)} \in \mathbb{R}^{d \times r}$ — a matrix whose columns are the top- r leading eigenvectors of $\mathbf{B}^{(m)}$. The procedure is summarized in Algorithm 5.

Algorithm 5 The m -th leave-one-out sequence for spectral initialization

- 1: Let $\mathbf{U}^{(m)} \mathbf{\Lambda}^{(m)} \mathbf{U}^{(m)\top}$ be the rank- r eigen-decomposition of $\mathbf{B}^{(m)}$ defined in (62).
 - 2: Generate the initial estimate $\mathbf{U}^{0,(m)} \in \mathbb{R}^{d \times r}$ from $\mathbf{U}^{(m)} \in \mathbb{R}^{d \times r}$ using Algorithm 6.
-

The following lemma plays a crucial role in our analysis, which formalizes the fact that the leave-one-out version $\mathbf{U}^{(m)}$ obtained by Algorithm 5 is extremely close to \mathbf{U} .

Lemma 5.8. *There exist some universal constants $c_0, c_1, c_2 > 0$ such that if*

$$p \geq c_0 \frac{\mu^2 r \log^2 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^*} \leq c_1 \frac{\sqrt{p}}{d^{3/4} \sqrt{\log d}}, \quad \text{and} \quad r \leq c_2 \sqrt{\frac{d}{\mu}},$$

then with probability $1 - O(d^{-10})$, the subspace estimate $\mathbf{U}^{(m)}$ computed by Algorithm 5 obeys

$$\|\mathbf{U} \mathbf{U}^\top - \mathbf{U}^{(m)} \mathbf{U}^{(m)\top}\|_{\text{F}} \lesssim \mathcal{E}_{\text{loo}} \sqrt{\frac{\mu r}{d}} \quad (64)$$

simultaneously for all $1 \leq m \leq d$, where

$$\mathcal{E}_{\text{loo}} := \frac{\mu^2 r \log d}{d^{3/2} p} + \sqrt{\frac{\mu^2 r \log d}{d^2 p}} + \frac{\sigma^2}{\lambda_{\min}^{*2}} \frac{d^{3/2} \log d}{p} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}}. \quad (65)$$

Lemma 5.8 follows immediately from the analysis of [CLC⁺20, Lemma 4]. As a remark, the construction of the leave-one-out sequences herein is slightly different from the one in [CLC⁺20]. However, it is straightforward to adapt the proof of [CLC⁺20] to the case considered herein. We therefore omit the proof for the sake of brevity.

5.3 Analysis for initialization: Part 2 (retrieval of individual tensor factors)

5.3.1 Main results and leave-one-out sequences

This section justifies that the procedure presented in Algorithm 3 allows to disentangle the tensor factors. For notational simplicity, we let

$$\bar{\mathbf{u}}_i^* := \mathbf{u}_i^* / \|\mathbf{u}_i^*\|_2, \quad \lambda_i^* := \|\mathbf{u}_i^*\|_2^3, \quad 1 \leq i \leq d. \quad (66)$$

Our result is this:

Theorem 5.9. Fix any arbitrary small constant $\delta > 0$. Assume that

$$p \geq c_0 \frac{\mu^2 r^4 \log^4 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^*} \leq c_1 \frac{\sqrt{p}}{r^{3/2} d^{3/4} \log^2 d}, \quad r \leq c_2 \left(\frac{d}{\mu^6 \log^6 d} \right)^{1/6},$$

$$L = c_3 r^{2\kappa^2} \log^{3/2} r, \quad \epsilon_{\text{th}} = c_4 \left\{ \frac{\mu r \log d}{d\sqrt{p}} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{rd \log^2 d}{p}} + \sqrt{\frac{\mu r \log d}{d}} \right\} \quad (67)$$

for some sufficiently large universal constant $c_0, c_3 > 0$ and some sufficiently small universal constants $c_1, c_2, c_4 > 0$. Then with probability exceeding $1 - \delta$, there exists a permutation $\pi(\cdot) : [d] \mapsto [d]$ such that for all $1 \leq i \leq r$, the tensor factors $\{\mathbf{w}^i\}_{i=1}^r$ returned by Algorithm 3 satisfy

$$\|\mathbf{w}^i - \bar{\mathbf{u}}_{\pi(i)}^*\|_2 \lesssim \frac{\mu r \log d}{d\sqrt{p}} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{rd \log^2 d}{p}} + \sqrt{\frac{\mu r \log d}{d}}; \quad (68a)$$

$$\|\mathbf{w}^i - \bar{\mathbf{u}}_{\pi(i)}^*\|_\infty \lesssim \left\{ \frac{\mu^2 r \log^4 d}{d^{3/2} p} + \frac{\mu r \log^3 d}{d\sqrt{p}} + \frac{\sigma^2}{\lambda_{\min}^{*2}} \frac{d^{3/2} \log^4 d}{p} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{rd \log^6 d}{p}} + \sqrt{\frac{\mu r \log^2 d}{d}} \right\} \sqrt{\frac{\mu r}{d}}; \quad (68b)$$

$$|\lambda_i - \lambda_{\pi(i)}^*| \lesssim \left\{ \frac{\mu r \log d}{d\sqrt{p}} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{rd \log^2 d}{p}} + \sqrt{\frac{\mu r \log d}{d}} \right\} \lambda_{\pi(i)}^*. \quad (68c)$$

In short, this theorem asserts that the estimates returned by Algorithm 3 are — up to global permutation — reasonably close to the ground truth under our sample size and noise conditions. In order to establish this theorem and in order to provide initial guesses for the leave-one-out GD sequences, we need to produce a leave-one-out sequence tailored to this part of the algorithm. Such auxiliary sequences are generated in a similar spirit as the previous ones, and we summarize them in Algorithm 6. As usual, the resulting leave-one-out estimates $\{\lambda_i^{(m)}, \mathbf{w}^{i,(m)}\}_{i=1}^r$ are statistically independent of $\mathcal{P}_{\Omega_m}(\mathbf{E})$.

In what follows, we gather a few key properties of the leave-one-out estimates, which play a crucial role in the analysis.

Algorithm 6 The m -th leave-one-out sequence for retrieving individual tensor components

- 1: **Input:** restart number L , threshold ϵ_{th} , subspace estimate $\mathbf{U}^{(m)} \in \mathbb{R}^{d \times r}$ given by Algorithm 5.
- 2: **for** $\tau = 1, \dots, L$ **do**
- 3: Recall the Gaussian vector $\mathbf{g}^\tau \sim \mathcal{N}(0, \mathbf{I}_d)$ generated in Algorithm 3.
- 4: $(\boldsymbol{\nu}^{\tau,(m)}, \lambda_\tau^{(m)}, \text{spec-gap}_\tau^{(m)}) \leftarrow \text{RETRIEVE-ONE-TENSOR-FACTOR}(\mathbf{T}^{(m)}, p, \mathbf{U}^{(m)}, \mathbf{g}^\tau)$.
- 5: Generate tensor factor estimates

$$\{(\mathbf{w}^{1,(m)}, \lambda_1^{(m)}), \dots, (\mathbf{w}^{r,(m)}, \lambda_r^{(m)})\} \leftarrow \text{PRUNE}(\{(\boldsymbol{\nu}^{\tau,(m)}, \lambda_\tau^{(m)}, \text{spec-gap}_\tau^{(m)})\}_{\tau=1}^L, \epsilon_{\text{th}}).$$

- 6: **Output:** an initial estimate $\mathbf{U}^{0,(m)} = [(\lambda_1^{(m)})^{1/3} \mathbf{w}^{1,(m)}, \dots, (\lambda_r^{(m)})^{1/3} \mathbf{w}^{r,(m)}]$.
-

Theorem 5.10. Fix any arbitrarily small constant $\delta > 0$. Instate the assumptions in Theorem 5.9. With probability exceeding $1 - \delta$, the permutation function stated in Theorem 5.9 obeys that: for all $1 \leq i \leq r$ and all $1 \leq m \leq d$:

$$\|\mathbf{w}^i - \mathbf{w}^{i,(m)}\|_2 \lesssim \left\{ \frac{\mu^2 r \log^{3/2} d}{d^{3/2} p} + \frac{\mu \sqrt{r} \log d}{d\sqrt{p}} + \frac{\sigma^2}{\lambda_{\min}^{*2}} \frac{d^{3/2} \log^{3/2} d}{p} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log^2 d}{p}} \right\} \sqrt{\frac{\mu r}{d}}; \quad (69a)$$

$$|\lambda_i - \lambda_i^{(m)}| \lesssim \left\{ \frac{\mu^2 r \log^{3/2} d}{d^{3/2} p} + \frac{\mu \sqrt{r} \log d}{d\sqrt{p}} + \frac{\sigma^2}{\lambda_{\min}^{*2}} \frac{d^{3/2} \log^{3/2} d}{p} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log^2 d}{p}} \right\} \sqrt{\frac{\mu r}{d}} \lambda_{\max}^*; \quad (69b)$$

$$|(\mathbf{w}^{i,(m)} - \bar{\mathbf{u}}_{\pi(i)}^*)| \lesssim \left\{ \frac{\sqrt{\mu r} \log^{7/2} d}{d^{3/2} p} + \frac{\mu r \log^3 d}{d \sqrt{p}} + \frac{\sigma \log^4 d}{\lambda_{\min}^* p} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{rd \log^6 d}{p}} + \sqrt{\frac{\mu r \log^2 d}{d}} \right\} \sqrt{\frac{\mu r}{d}}. \quad (69c)$$

With Theorems 5.9-5.10 in place, we can immediately establish a few desired properties (particularly those specified in Section 5.1) of our initial estimate, as asserted in the following corollary.

Corollary 5.11. *Fix any arbitrarily small constant $\delta > 0$. Instate the assumptions in Theorem 2.8. With probability exceeding $1 - \delta$, the estimates \mathbf{U}^0 and $\mathbf{U}^{0,(m)}$ returned by Algorithm 3 and Algorithm 6 respectively satisfy the hypotheses (47) for $t = 0$.*

Proof. See Appendix B.10. ■

5.3.2 Analysis

Before we start with the proof, we first state the main idea. For the sake of clarify, we define

$$\boldsymbol{\theta}^\tau := \mathbf{U} \mathbf{U}^\top \mathbf{g}^\tau, \quad (70a)$$

$$\boldsymbol{\theta}^{\tau,(m)} := \mathbf{U}^{(m)} \mathbf{U}^{(m)\top} \mathbf{g}^\tau, \quad (70b)$$

$$\mathbf{M}^\tau := p^{-1} \mathbf{T} \times_3 \boldsymbol{\theta}^\tau, \quad (70c)$$

$$\mathbf{M}^{\tau,(m)} := p^{-1} \mathbf{T}^{(m)} \times_3 \boldsymbol{\theta}^{\tau,(m)}. \quad (70d)$$

In addition, let $\boldsymbol{\nu}^\tau$ be the top singular vector of \mathbf{M}^τ obeying $\langle \mathbf{T}, (\boldsymbol{\nu}^\tau)^{\otimes 3} \rangle \geq 0$, and $\boldsymbol{\nu}^{\tau,(m)}$ the top singular vector of $\mathbf{M}^{\tau,(m)}$ obeying $\langle \mathbf{T}^{(m)}, (\boldsymbol{\nu}^{\tau,(m)})^{\otimes 3} \rangle \geq 0$. Set

$$\lambda_\tau := \langle p^{-1} \mathbf{T}, (\boldsymbol{\nu}^\tau)^{\otimes 3} \rangle \quad \text{and} \quad \lambda_\tau^{(m)} := \langle p^{-1} \mathbf{T}^{(m)}, (\boldsymbol{\nu}^{\tau,(m)})^{\otimes 3} \rangle. \quad (71)$$

These are all computed in the function RETRIEVE-ONE-TENSOR-FACTOR() in the τ -th round.

1. We first show that for each $1 \leq i \leq r$, there exists at least one trial $1 \leq \tau \leq L$ such that the i -th tensor factor $\bar{\mathbf{u}}_i^*$ is the top singular vector of the population version of $\mathbf{T} \times_3 \boldsymbol{\theta}^\tau$ (with respect to the missing data and noise). In addition, the spectral gap is large enough to guarantee accurate estimates.
2. Next, we prove that given this spectral gap, the top singular vector $\boldsymbol{\nu}^\tau$ of $\mathbf{T} \times_3 \boldsymbol{\theta}^\tau$ is close to $\bar{\mathbf{u}}_i^*$ both in the $\|\cdot\|_2$ and $\|\cdot\|_\infty$ norm. This also enables us to accurately estimate the magnitude of \mathbf{u}_i^* .
3. Finally, we need to show that one can find those reliable estimates among L random restarts. Combining the spectral gap information with the incoherence condition that tensor components are nearly orthogonal to each other, our selection procedure is guaranteed to recover all tensor factors.

Now we proceed to the proof. Without loss of generality, we prove the case for $i = 1$ in the sequel, i.e. there exists some $\tau \in [L]$ such that $\boldsymbol{\nu}^\tau$ accurately recovers \mathbf{u}_1^* . Together with the union bound, this shows that we can find reliable estimates for all tensor factors. We then conclude the proof by showing that Algorithm 3 is able to find all of them without duplicates.

To this end, we find it convenient to introduce an auxiliary vector $\boldsymbol{\gamma}^{*\tau} = [\gamma_1^{*\tau}, \dots, \gamma_r^{*\tau}]^\top \in \mathbb{R}^r$ and its leave-one-out versions $\boldsymbol{\gamma}^{*\tau,(m)} = [\gamma_1^{*\tau,(m)}, \dots, \gamma_r^{*\tau,(m)}]^\top$ ($1 \leq m \leq d$) for each $1 \leq \tau \leq L$ as follows:

$$\gamma_i^{*\tau} := \|\mathbf{u}_i^*\|_2^2 \langle \mathbf{u}_i^*, \boldsymbol{\theta}^\tau \rangle = \lambda_i^* \langle \bar{\mathbf{u}}_i^*, \boldsymbol{\theta}^\tau \rangle, \quad (72a)$$

$$\gamma_i^{*\tau,(m)} := \|\mathbf{u}_i^*\|_2^2 \langle \mathbf{u}_i^*, \boldsymbol{\theta}^{\tau,(m)} \rangle = \lambda_i^* \langle \bar{\mathbf{u}}_i^*, \boldsymbol{\theta}^{\tau,(m)} \rangle, \quad (72b)$$

where $\bar{\mathbf{u}}_i^*$ and λ_i^* are both defined in (66). The idea is to let $\boldsymbol{\gamma}^{*\tau}$ approximate the singular values of $\mathbf{T}^* \times_3 \boldsymbol{\theta}^\tau$; this can be seen, for instance, via the following calculation:

$$\mathbf{T}^* \times_3 \boldsymbol{\theta}^\tau = \sum_{i=1}^r \langle \mathbf{u}_i^*, \boldsymbol{\theta}^\tau \rangle \mathbf{u}_i^* \mathbf{u}_i^{*\top} = \sum_{i=1}^r \underbrace{\lambda_i^* \langle \bar{\mathbf{u}}_i^*, \boldsymbol{\theta}^\tau \rangle}_{=\gamma_i^{*\tau}} \bar{\mathbf{u}}_i^* \bar{\mathbf{u}}_i^{*\top}, \quad (73)$$

where $\{\bar{\mathbf{u}}_i^*\}_{i=1}^r$ — which are assumed to be incoherent (or nearly orthogonal to each other) — can be approximately viewed as the singular vectors of $\mathbf{T}^* \times_3 \boldsymbol{\theta}^\tau$.

If we want our spectral estimate to be accurate, we would need to be assured that the two largest entries of $\boldsymbol{\gamma}^{*\tau}$ (in magnitude) are sufficiently separated.

Lemma 5.12. *Instate the assumptions of Theorem 5.9. Define $\Delta_1^\tau := \gamma_1^{*\tau} - \max_{1 < i \leq r} |\gamma_i^{*\tau}|$ for each $1 \leq \tau \leq L$ and let $\Delta_1^{(1)} \geq \Delta_1^{(2)} \geq \dots \geq \Delta_1^{(L)}$ denote the order statistics of $\{\Delta_1^\tau\}_{\tau=1}^L$ (in descending order). Fix any arbitrary small constant $\delta > 0$. With probability greater than $1 - \delta/r$, one has*

$$\Delta_1^{(1)} \gtrsim \lambda_{\min}^*, \quad (74a)$$

$$\Delta_1^{(1)} - \Delta_1^{(2)} \gtrsim \frac{\lambda_{\min}^*}{r\sqrt{\log d}}, \quad (74b)$$

Additionally, for any fixed vector $\mathbf{v} \in \mathbb{R}^r$, with probability at least $1 - O(d^{-10})$, for all $1 \leq \tau \leq L$, one has

$$\gamma_1^{*\tau} \lesssim \sqrt{\log d} \lambda_{\max}^*, \quad (75a)$$

$$\|\boldsymbol{\gamma}^{*\tau}\|_2 \lesssim \sqrt{r \log d} \lambda_{\max}^*, \quad (75b)$$

$$|\langle \mathbf{v}, \boldsymbol{\gamma}^{*\tau} \rangle| \lesssim \|\mathbf{v}\|_2 \sqrt{\log d} \lambda_{\max}^*. \quad (75c)$$

Proof. See Appendix B.1. ■

Lemma 5.12 demonstrates that there exists some $\tau \in [L]$ such that $\gamma_1^{*\tau} - \max_{1 < i \leq r} |\gamma_i^{*\tau}| \gtrsim \lambda_{\min}^*$. This means that $\bar{\mathbf{u}}_1^*$ exhibits the largest correlation with the random projection $\boldsymbol{\theta}$, which further implies that $\bar{\mathbf{u}}_1^*$ is the largest singular vector of $\mathbf{T}^* \times_3 \boldsymbol{\theta}^\tau$ with a considerable spectral gap (as we will show shortly). With the desired spectral gap in place, we are ready to look at the eigenvectors / singular vectors of interest. To this end, we find it convenient to introduce another auxiliary vector $\bar{\mathbf{u}}^\tau$, defined as the leading singular vector of \mathbf{M}^τ (cf. (70e)) obeying

$$\langle \bar{\mathbf{u}}^\tau, \bar{\mathbf{u}}_1^* \rangle \geq 0. \quad (76)$$

The careful reader would immediately notice the similarity between $\bar{\mathbf{u}}^\tau$ and $\boldsymbol{\nu}^\tau$ except for their global signs; namely, we determine the global sign of $\bar{\mathbf{u}}^\tau$ based on the ground truth information (76), but pick the global sign for $\boldsymbol{\nu}^\tau$ solely based on the observed data (cf. Algorithm 3). Fortunately, the vectors $\bar{\mathbf{u}}^\tau$ and $\boldsymbol{\nu}^\tau$ provably coincide, namely,

$$\bar{\mathbf{u}}^\tau = \boldsymbol{\nu}^\tau, \quad (77)$$

as we shall demonstrate momentarily in Lemma 5.16. In a similar way, we also denote by $\bar{\mathbf{u}}^{\tau,(m)}$ the leading singular vector of $\mathbf{M}^{\tau,(m)}$ defined in (70d) such that

$$\langle \bar{\mathbf{u}}^{\tau,(m)}, \bar{\mathbf{u}}_1^* \rangle \geq 0. \quad (78)$$

Lemma 5.16 also shows that $\bar{\mathbf{u}}^{\tau,(m)} = \boldsymbol{\nu}^{\tau,(m)}$.

We shall now take a detour to look at $\bar{\mathbf{u}}^\tau$, which in turn would help us understand $\boldsymbol{\nu}^\tau$. We shall first demonstrate that $\bar{\mathbf{u}}^\tau$ (and hence $\boldsymbol{\nu}^\tau$) is sufficiently close to the corresponding true factor in the ℓ_2 sense.

Lemma 5.13. *Instate the assumptions of Theorem 5.9. Let $\bar{\mathbf{u}}^\tau$ and $\bar{\mathbf{u}}_1^*$ be as defined in (76) and (66), respectively. Define \mathcal{A} to be the event such that $\gamma_1^{*\tau} - \max_{1 < i \leq r} |\gamma_i^{*\tau}| \gtrsim \lambda_{\min}^*$ and the condition (75) hold. Then conditional on this event \mathcal{A} , with probability exceeding $1 - O(d^{-11})$ one has*

$$\|\bar{\mathbf{u}}^\tau - \bar{\mathbf{u}}_1^*\|_2 \lesssim \underbrace{\frac{\mu r \log d}{d\sqrt{p}} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{rd \log^2 d}{p}} + \sqrt{\frac{\mu r \log d}{d}}}_{=: \mathcal{E}_{\text{proj}}}. \quad (79)$$

Proof. See Appendix B.2. ■

Thus far, we have focused on the ℓ_2 estimation errors. In order to further quantify the ℓ_∞ estimation errors, we need to resort to the leave-one-out estimates $\bar{\mathbf{u}}^{\tau,(m)}$ ($1 \leq m \leq d$). Specifically, we shall justify in the following two lemmas that: (1) the m -th leave-one-out estimate $\bar{\mathbf{u}}^{\tau,(m)}$ is close to the truth at least in the m -th coordinate, and (2) the vector $\bar{\mathbf{u}}^\tau$ is extremely close to each of the leave-one-out estimates $\bar{\mathbf{u}}^{\tau,(m)}$ ($1 \leq m \leq d$). These two observations taken collectively translate to the desired entrywise error control of $\bar{\mathbf{u}}^\tau$. Here, we recall that the global sign of $\bar{\mathbf{u}}^{\tau,(m)}$ (cf. (78)) and the global sign of $\bar{\mathbf{u}}^\tau$ (cf. (76)) are defined in a similar fashion, both using the ground truth information.

Lemma 5.14. *Instate the assumptions of Theorem 5.9. Define \mathcal{A} to be the event such that $\gamma_1^{*\tau} - \max_{1 < i \leq r} |\gamma_i^{*\tau}| \gtrsim \lambda_{\min}^*$ and the condition (75) hold. Then conditional on this event \mathcal{A} , one has, with probability exceeding $1 - O(d^{-10})$, that*

$$\left| [\bar{\mathbf{u}}^{\tau,(m)} - \bar{\mathbf{u}}_1^*]_m \right| \lesssim \mathcal{E}_{\text{op}} \sqrt{\frac{\mu r \log d}{d}}, \quad (80)$$

holds for all $m \in [d]$, where \mathcal{E}_{op} is defined as follows:

$$\mathcal{E}_{\text{op}} := \frac{\sqrt{\mu r} \log^3 d}{d^{3/2} p} + \frac{\mu r \log^{5/2} d}{d \sqrt{p}} + \frac{\sigma \log^{7/2} d}{\lambda_{\min}^* p} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{rd \log^5 d}{p}} + \sqrt{\frac{\mu r \log d}{d}}. \quad (81)$$

Proof. See Appendix B.4. ■

Lemma 5.15. *Instate the assumptions of Theorem 5.9. Define \mathcal{A} to be the event such that $\gamma_1^{*\tau} - \max_{1 < i \leq r} |\gamma_i^{*\tau}| \gtrsim \lambda_{\min}^*$ and the condition (75) hold. Then conditional on this event \mathcal{A} , one has, with probability at least $1 - O(d^{-10})$, for all $m \in [d]$:*

$$\|\bar{\mathbf{u}}^\tau - \bar{\mathbf{u}}^{\tau,(m)}\|_2 \lesssim \mathcal{E}_{\text{loo}} \sqrt{\frac{\mu r \log d}{d}}, \quad (82)$$

$$\|\bar{\mathbf{u}}^\tau - \bar{\mathbf{u}}_1^*\|_\infty \lesssim (\mathcal{E}_{\text{op}} + \mathcal{E}_{\text{loo}}) \sqrt{\frac{\mu r \log d}{d}}, \quad (83)$$

$$|\lambda_\tau - \lambda_\tau^{(m)}| \lesssim \mathcal{E}_{\text{loo}} \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^*, \quad (84)$$

where \mathcal{E}_{op} and \mathcal{E}_{loo} are defined in (81) and (65), respectively.

Proof. See Appendix B.6. ■

Next, we turn to the estimation accuracy regarding the size of the tensor factors and show that λ_τ (produced in Algorithm 3) is close to the truth as well. As it turns out, a byproduct of this step reveals that $\boldsymbol{\nu}^\tau = \bar{\mathbf{u}}^\tau$ and $\boldsymbol{\nu}^{\tau,(m)} = \bar{\mathbf{u}}^{\tau,(m)}$, where $\bar{\mathbf{u}}^\tau$ and $\bar{\mathbf{u}}^{\tau,(m)}$ are an auxiliary vectors defined in (76) and (78), respectively.

Lemma 5.16. *Instate the assumptions of Theorem 5.9. Assume that the results in Lemma 5.13, Lemma 5.14 and Lemma 5.15 hold. Then with probability at least $1 - O(d^{-10})$, one has*

$$|\lambda_\tau - \lambda_1^*| \lesssim \mathcal{E}_{\text{proj}} \lambda_1^*. \quad (85)$$

In particular, one has

$$\boldsymbol{\nu}^\tau = \bar{\mathbf{u}}^\tau \quad \text{and} \quad \boldsymbol{\nu}^{\tau,(m)} = \bar{\mathbf{u}}^{\tau,(m)}, \quad 1 \leq m \leq d. \quad (86)$$

Proof. See Appendix B.8. ■

Thus far, we have only proved that one can find a reliable estimate for each tensor factor within L random trials, provided that L is sufficiently large. To finish up, it remains to show that the pruning procedure PRUNE() is capable of returning a rough estimate for each tensor factor without duplication. This is accomplished in the following lemma.

Lemma 5.17. *Instate the assumptions of Theorem 5.9. On the event that the results in Lemma 5.13, Lemma 5.14, Lemma 5.15 and Lemma 5.16 hold for all $1 \leq i \leq r$, there exists a permutation $\pi(\cdot) : [d] \mapsto [d]$ such that: for each $1 \leq i \leq r$, $(\lambda_i, \mathbf{w}^i)$ and $(\lambda_{\pi(i)}^*, \bar{\mathbf{u}}_{\pi(i)}^*)$ satisfy (68a), (68b) and (68c); $(\lambda_i, \mathbf{w}^i)$ and $(\lambda_i^{(m)}, \mathbf{w}^{i,(m)})_{i=1}^r$ obey (69a), (69b) and (69c) for all $1 \leq m \leq d$, where $\{\lambda_i, \mathbf{w}^i\}_{i=1}^r$ and $\{\lambda_i^{(m)}, \mathbf{w}^{i,(m)}\}_{i=1}^r$ are outputs of Algorithm 3 and Algorithm 6, respectively.*

Proof. See Appendix B.9. ■

6 Discussion

The current paper uncovers the possibility of efficiently and stably completing a low-CP-rank tensor from partial and noisy entries. Perhaps somewhat unexpectedly, despite the high degree of nonconvexity, this problem can be solved to optimal statistical accuracy within nearly linear time, provided that the tensor of interest is well-conditioned, incoherent, and of constant rank. To the best of our knowledge, this intriguing message has not been shown in the prior literature.

Moving forward, one pressing issue is to understand how to improve the algorithmic and theoretical dependency upon the tensor rank r of the proposed method. Ideally one would desire a fast algorithm whose sample complexity scales as $rd^{1.5}$, an order that is provably achievable by the sum-of-squares hierarchy. Additionally, in contrast to the matrix counterpart where the rank is upper bounded by the matrix dimension, the tensor CP rank is allowed to rise above d , which is commonly referred to as the over-complete case. Unfortunately, our current initialization scheme (i.e. the spectral method) fails to work unless $r < d$, and our local analysis for GD falls of accommodating the scenario with $r > d$. It would be of great interest to develop more powerful algorithms — in addition to more refined analysis — to tackle such an important over-complete regime.

Another tantalizing research direction is the exploration of landscape design for tensor completion. As our heuristic discussions as well as other prior work (e.g. [RM14]) suggest, randomly initialized gradient descent tailored to (4) seems unlikely to work, unless the sample size is significantly larger than the computational limit. This might mean either that there exist spurious local minima in the natural nonconvex least squares formulation (4), or that the optimization landscape of (4) is too flat around some saddle points and hence not amenable to fast computation. It would be interesting to investigate what families of loss functions allow us to rule out bad local minima and eliminate the need of careful initialization, which might be better suited for tensor recovery problems.

Finally, in statistical inference and decision making, one might not be simply satisfied with obtaining a reliable estimate for each missing entry, but would also like to report a short confidence interval which is likely to contain the true entry. This boils down to the fundamental task of uncertainty quantification for tensor completion, which we leave to future investigation.

Acknowledgements

Y. Chen is supported in part by the AFOSR YIP award FA9550-19-1-0030, by the ONR grant N00014-19-1-2120, by the ARO grants W911NF-20-1-0097 and W911NF-18-1-0303, by the NSF grants CCF-1907661, IIS-1900140 and DMS-2014279, and by the Princeton SEAS innovation award. H. V. Poor is supported in part by the NSF grant DMS-1736417. C. Cai is supported in part by Gordon Y. S. Wu Fellowships in Engineering. This work was done in part while Y. Chen was visiting the Kavli Institute for Theoretical Physics (supported in part by NSF grant PHY-1748958). We thank Lanqing Yu for many helpful discussions, and thank Yuling Yan for proofreading the paper.

A Proofs for local convergence of GD

In this section, we establish the key lemmas concerning the convergence properties of GD. As one can easily see, treating $\{E_{i,j,k}\}_{1 \leq i,j,k \leq d}$ (resp. $\{\chi_{i,j,k}\}_{1 \leq i,j,k \leq d}$) as independent random variables — which leads to

asymmetric versions of \mathbf{E} and Ω — does not affect the order of our results at all. In light of this, we shall adopt such an independent assumption whenever it simplifies our presentation.

A.1 Proof of Lemma 5.1

For notational convenience, for any matrix $\mathbf{M} = [\mathbf{m}_1, \dots, \mathbf{m}_r] \in \mathbb{R}^{d \times r}$, let

$$\widetilde{\mathbf{M}} := [\mathbf{m}_1 \otimes \mathbf{m}_1, \dots, \mathbf{m}_r \otimes \mathbf{m}_r] \in \mathbb{R}^{d^2 \times r}, \quad (87)$$

where for any $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ we denote $\mathbf{a} \otimes \mathbf{b} := \begin{bmatrix} a_1 \mathbf{b} \\ \vdots \\ a_d \mathbf{b} \end{bmatrix} \in \mathbb{R}^{d^2}$.

From the Hessian expression (41), one can decompose

$$\begin{aligned} & \text{vec}(\mathbf{V})^\top \nabla^2 f(\mathbf{U}) \text{vec}(\mathbf{V}) \\ &= \underbrace{\frac{1}{3p} \left\| \mathcal{P}_\Omega \left(\sum_{s \in [r]} (\mathbf{v}_s \otimes \mathbf{u}_s^{\otimes 2} + \mathbf{u}_s \otimes \mathbf{v}_s \otimes \mathbf{u}_s + \mathbf{u}_s^{\otimes 2} \otimes \mathbf{v}_s) \right) \right\|_{\mathbb{F}}^2 - \frac{1}{3p} \left\| \mathcal{P}_\Omega \left(\sum_{s \in [r]} (\mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} + \mathbf{u}_s^* \otimes \mathbf{v}_s \otimes \mathbf{u}_s^* + \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{v}_s) \right) \right\|_{\mathbb{F}}^2}_{=: \alpha_1}} \\ &+ \underbrace{\frac{1}{3p} \left\| \mathcal{P}_\Omega \left(\sum_{s \in [r]} (\mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} + \mathbf{u}_s^* \otimes \mathbf{v}_s \otimes \mathbf{u}_s^* + \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{v}_s) \right) \right\|_{\mathbb{F}}^2 - \frac{1}{3} \left\| \sum_{s \in [r]} (\mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} + \mathbf{u}_s^* \otimes \mathbf{v}_s \otimes \mathbf{u}_s^* + \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{v}_s) \right\|_{\mathbb{F}}^2}_{=: \alpha_2}} \\ &+ 2 \underbrace{\left\langle \frac{1}{p} \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{u}_s^{\otimes 3} - \mathbf{T}^* \right), \sum_{s \in [r]} \mathbf{v}_s^{\otimes 2} \otimes \mathbf{u}_s \right\rangle}_{=: \alpha_3} + \underbrace{\frac{1}{3} \left\| \sum_{s \in [r]} (\mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} + \mathbf{u}_s^* \otimes \mathbf{v}_s \otimes \mathbf{u}_s^* + \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{v}_s) \right\|_{\mathbb{F}}^2}_{=: \alpha_4}. \end{aligned}$$

In what follows, we shall bound each of the above terms separately.

A.1.1 Bounding α_4

With regards to α_4 , by symmetry we have

$$\alpha_4 = \left\| \sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} \right\|_{\mathbb{F}}^2 + 2 \left\langle \sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2}, \sum_{s \in [r]} \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{v}_s \right\rangle. \quad (88)$$

In order to control (88), we first see that

$$\left\| \sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} \right\|_{\mathbb{F}} = \|\mathbf{V} \widetilde{\mathbf{U}}^{*\top}\|_{\mathbb{F}}, \quad (89)$$

where $\widetilde{\mathbf{U}}^*$ is as defined in (87). Similar to the proof of Lemma D.1, we can use the fact that $\langle \mathbf{u}_i^{*\otimes 2}, \mathbf{u}_j^{*\otimes 2} \rangle = \langle \mathbf{u}_i^*, \mathbf{u}_j^* \rangle^2$ and (8c) to deduce that

$$\sigma_{\min}(\widetilde{\mathbf{U}}^*) = \lambda_{\min}^{*2/3} (1 + o(1)) \quad \text{and} \quad \sigma_{\max}(\widetilde{\mathbf{U}}^*) = \lambda_{\max}^{*2/3} (1 + o(1)), \quad (90)$$

provided that $r \ll d/\mu$. This implies that

$$\frac{19}{20} \lambda_{\min}^{*2/3} \|\mathbf{V}\|_{\mathbb{F}} \leq \sigma_{\min}(\widetilde{\mathbf{U}}^*) \|\mathbf{V}\|_{\mathbb{F}} \leq \|\mathbf{V} \widetilde{\mathbf{U}}^{*\top}\|_{\mathbb{F}} \leq \sigma_{\max}(\widetilde{\mathbf{U}}^*) \|\mathbf{V}\|_{\mathbb{F}} \leq \frac{11}{10} \lambda_{\max}^{*2/3} \|\mathbf{V}\|_{\mathbb{F}}. \quad (91)$$

(1) Speaking of an upper bound on α_4 , we can invoke the Cauchy-Schwarz inequality followed by (91) to reach

$$\alpha_4 \leq 3 \left\| \sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} \right\|_{\mathbb{F}}^2 = 3 \|\mathbf{V} \widetilde{\mathbf{U}}^{*\top}\|_{\mathbb{F}}^2 \leq \frac{7}{2} \lambda_{\max}^{*4/3} \|\mathbf{V}\|_{\mathbb{F}}^2. \quad (92)$$

(2) When it comes to lower bounding α_4 , the main step boils down to controlling the inner product term in (88). Applying the Cauchy-Schwartz inequality gives that

$$\begin{aligned}
\left\langle \sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2}, \sum_{s \in [r]} \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{v}_s \right\rangle &= \sum_{s \in [r]} \langle \mathbf{v}_s, \mathbf{u}_s^* \rangle^2 \|\mathbf{u}_s^*\|_2^2 + \sum_{s_1 \neq s_2} \langle \mathbf{v}_{s_1}, \mathbf{u}_{s_2}^* \rangle \langle \mathbf{u}_{s_1}^*, \mathbf{v}_{s_2} \rangle \langle \mathbf{u}_{s_1}^*, \mathbf{u}_{s_2}^* \rangle \\
&\geq - \max_{s_1 \neq s_2} |\langle \mathbf{u}_{s_1}^*, \mathbf{u}_{s_2}^* \rangle| \sum_{s_1 \neq s_2} \|\mathbf{v}_{s_1}\|_2 \|\mathbf{u}_{s_1}^*\|_2 \|\mathbf{v}_{s_2}\|_2 \|\mathbf{u}_{s_2}^*\|_2 \\
&\geq - \max_{s_1 \neq s_2} |\langle \mathbf{u}_{s_1}^*, \mathbf{u}_{s_2}^* \rangle| \left(\sum_{s \in [r]} \|\mathbf{v}_s\|_2 \|\mathbf{u}_s^*\|_2 \right)^2 \\
&\stackrel{(i)}{\geq} - \max_{s_1 \neq s_2} |\langle \mathbf{u}_{s_1}^*, \mathbf{u}_{s_2}^* \rangle| \|\mathbf{U}^*\|_{\text{F}}^2 \|\mathbf{V}\|_{\text{F}}^2 \\
&\geq -r \sqrt{\frac{\mu}{d}} \lambda_{\max}^{*4/3} \|\mathbf{V}\|_{\text{F}}^2 \geq -\frac{1}{40} \lambda_{\min}^{*4/3} \|\mathbf{V}\|_{\text{F}}^2,
\end{aligned}$$

where (i) comes from Cauchy-Schwartz, and the last line follows from (8c), (8d) as well as the condition that $r \ll \sqrt{d/\mu}$ and $\kappa \asymp 1$. Therefore, we can lower bound α_4 by (with the assistance of (91))

$$\begin{aligned}
\alpha_4 &= \left\| \sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} \right\|_{\text{F}}^2 + 2 \left\langle \sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2}, \sum_{s \in [r]} \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{v}_s \right\rangle \\
&\geq \frac{19}{20} \lambda_{\min}^{*4/3} \|\mathbf{V}\|_{\text{F}}^2 - \frac{1}{20} \lambda_{\min}^{*4/3} \|\mathbf{V}\|_{\text{F}}^2 \geq \frac{9}{10} \lambda_{\min}^{*4/3} \|\mathbf{V}\|_{\text{F}}^2.
\end{aligned} \tag{93}$$

A.1.2 Bounding α_1

When it comes to α_1 , we can expand

$$\begin{aligned}
&\left\| \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{\otimes 2} \right) + \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{u}_s \otimes \mathbf{v}_s \otimes \mathbf{u}_s \right) + \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{u}_s^{\otimes 2} \otimes \mathbf{v}_s \right) \right\|_{\text{F}}^2 \\
&= 3 \left\| \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{\otimes 2} \right) \right\|_{\text{F}}^2 + 6 \left\langle \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{\otimes 2} \right), \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{u}_s^{\otimes 2} \otimes \mathbf{v}_s \right) \right\rangle;
\end{aligned}$$

we can decompose $\left\| \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} \right) + \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{u}_s^* \otimes \mathbf{v}_s \otimes \mathbf{u}_s^* \right) + \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{v}_s \right) \right\|_{\text{F}}^2$ in a similar way. As a consequence,

$$\begin{aligned}
\alpha_1 &= \frac{2}{p} \left(\underbrace{\left\langle \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{\otimes 2} \right), \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{u}_s^{\otimes 2} \otimes \mathbf{v}_s \right) \right\rangle - \left\langle \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} \right), \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{v}_s \right) \right\rangle}_{=:\beta_1} \right) \\
&\quad + \frac{1}{p} \left(\underbrace{\left\| \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{\otimes 2} \right) \right\|_{\text{F}}^2 - \left\| \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} \right) \right\|_{\text{F}}^2}_{=:\beta_2} \right).
\end{aligned}$$

We will derive an upper bound on β_1 in the sequel; the same method immediately applies to β_2 .

For notational convenience, let us define

$$\mathbf{\Delta} := \mathbf{U} - \mathbf{U}^*, \quad \mathbf{\Delta}_s := \mathbf{u}_s - \mathbf{u}_s^*, \quad \tilde{\mathbf{\Delta}} := [\mathbf{\Delta}_1 \otimes \mathbf{\Delta}_1, \dots, \mathbf{\Delta}_r \otimes \mathbf{\Delta}_r] \in \mathbb{R}^{d^2 \times r}. \tag{94}$$

Then one can write

$$\begin{aligned}
\frac{1}{2} \beta_1 &= \frac{1}{p} \left\langle \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{v}_s \otimes (\mathbf{\Delta}_s + \mathbf{u}_s^*)^{\otimes 2} \right), \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} (\mathbf{\Delta}_s + \mathbf{u}_s^*)^{\otimes 2} \otimes \mathbf{v}_s \right) \right\rangle - \frac{1}{p} \left\langle \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} \right), \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{v}_s \right) \right\rangle \\
&= \frac{1}{p} \left\langle \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} \right), \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{u}_s^* \otimes \mathbf{\Delta}_s \otimes \mathbf{v}_s \right) + \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{\Delta}_s \otimes \mathbf{u}_s^* \otimes \mathbf{v}_s \right) + \mathcal{P}_{\Omega} \left(\sum_{s \in [r]} \mathbf{\Delta}_s \otimes \mathbf{\Delta}_s \otimes \mathbf{v}_s \right) \right\rangle
\end{aligned}$$

$$\begin{aligned}
& + \frac{1}{p} \left\langle \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^* \otimes \Delta_s \right) + \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \Delta_s \otimes \mathbf{u}_s^* \right) + \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \Delta_s^{\otimes 2} \right), \right. \\
& \left. \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{v}_s \right) + \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{u}_s^* \otimes \Delta_s \otimes \mathbf{v}_s \right) + \mathcal{P}_\Omega \left(\sum_{s \in [r]} \Delta_s \otimes \mathbf{u}_s^* \otimes \mathbf{v}_s \right) + \mathcal{P}_\Omega \left(\sum_{s \in [r]} \Delta_s \otimes \Delta_s \otimes \mathbf{v}_s \right) \right\rangle.
\end{aligned}$$

Apply the Cauchy-Schwartz inequality to yield that

$$\begin{aligned}
|\beta_1| & \lesssim \frac{1}{p} \left\| \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} \right) \right\|_{\text{F}} \left(2 \left\| \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{u}_s^* \otimes \Delta_s \otimes \mathbf{v}_s \right) \right\|_{\text{F}} + \left\| \mathcal{P}_\Omega \left(\sum_{s \in [r]} \Delta_s \otimes \Delta_s \otimes \mathbf{v}_s \right) \right\|_{\text{F}} \right) \\
& + \frac{1}{p} \left(\left\| \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{v}_s \right) \right\|_{\text{F}} + 2 \left\| \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{u}_s^* \otimes \Delta_s \otimes \mathbf{v}_s \right) \right\|_{\text{F}} + \left\| \mathcal{P}_\Omega \left(\sum_{s \in [r]} \Delta_s \otimes \Delta_s \otimes \mathbf{v}_s \right) \right\|_{\text{F}} \right) \\
& \cdot \left(2 \left\| \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^* \otimes \Delta_s \right) \right\|_{\text{F}} + \left\| \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \Delta_s^{\otimes 2} \right) \right\|_{\text{F}} \right). \tag{95}
\end{aligned}$$

Before we bound the above quantities, we pause to make the following observations. In view of the assumptions of this lemma that $\delta \ll 1/\sqrt{r} \leq 1$, the following holds for all $i \in [r]$:

$$\|\Delta_i\|_2 \leq \|\mathbf{U} - \mathbf{U}^*\|_{\text{F}} \leq \delta \|\mathbf{U}^*\|_{\text{F}} \leq \delta \sqrt{r} \lambda_{\max}^{*1/3} \ll \lambda_{\max}^{*1/3}, \tag{96a}$$

$$\|\Delta_i\|_\infty \leq \|\mathbf{U} - \mathbf{U}^*\|_{2,\infty} \leq \delta \|\mathbf{U}^*\|_{2,\infty} \leq \delta \sqrt{\frac{\mu r}{d}} \lambda_{\max}^{*1/3} \ll \sqrt{\frac{\mu}{d}} \lambda_{\max}^{*1/3}, \tag{96b}$$

$$\|\mathbf{u}_i\|_2 \leq \|\mathbf{u}_i^*\|_2 + \|\Delta_i\|_2 \leq 2\lambda_{\max}^{*1/3}, \tag{96c}$$

$$\|\mathbf{u}_i\|_\infty \leq \|\mathbf{u}_i^*\|_\infty + \|\Delta_i\|_\infty \leq 2\sqrt{\frac{\mu}{d}} \lambda_{\max}^{*1/3}, \tag{96d}$$

Consequently, we also know that

$$\|\tilde{\Delta}\|_{2,\infty} \leq \max_{1 \leq i \leq r} \|\Delta_i\|_\infty \|\Delta\|_{2,\infty} \leq \delta^2 \|\mathbf{U}^*\|_{2,\infty}^2. \tag{97}$$

Now, we proceed to prove the claim. Let us define $\mathcal{S}_i := \{j \in [d]^2 \mid \chi_{ij_1 j_2} = 1\}$ for each $i \in [d]$. Applying the Chernoff bound and the union bound yields that: with probability at least $1 - O(d^{-10})$ one has

$$\max_{i \in [d]} |\mathcal{S}_i| \lesssim d^2 p. \tag{98}$$

provided $p \gg d^{-2} \log d$. It then follows from the Cauchy-Schwarz inequality that

$$\begin{aligned}
\frac{1}{p} \left\| \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \Delta_s^{\otimes 2} \right) \right\|_{\text{F}}^2 & = \frac{1}{p} \sum_{i \in [d], \underline{j} \in [d]^2} \chi_{ij_1 j_2} \langle \mathbf{v}_{i,:}, \tilde{\Delta}_{\underline{j},:} \rangle^2 \\
& \leq \frac{1}{p} \sum_{i \in [d]} \|\mathbf{v}_{i,:}\|_2^2 \sum_{\underline{j} \in \mathcal{S}_i} \|\tilde{\Delta}_{\underline{j},:}\|_2^2 \\
& \leq \frac{1}{p} \max_{i \in [d]} |\mathcal{S}_i| \|\tilde{\Delta}\|_{2,\infty}^2 \|\mathbf{V}\|_{\text{F}}^2 \\
& \stackrel{(i)}{\lesssim} d^2 \delta^4 \|\mathbf{U}^*\|_{2,\infty}^4 \|\mathbf{V}\|_{\text{F}}^2 \lesssim \delta^4 \mu^2 r^2 \lambda_{\max}^{*4/3} \|\mathbf{V}\|_{\text{F}}^2,
\end{aligned}$$

where (i) arises from (97) and (98). In a similar manner, we can derive

$$\begin{aligned}
\frac{1}{p} \left\| \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{u}_s^* \otimes \mathbf{v}_s \otimes \Delta_s \right) \right\|_{\text{F}}^2 & \lesssim \frac{1}{p} \max_{1 \leq i \leq d} |\mathcal{S}_i| \max_{1 \leq s \leq r} \|\mathbf{u}_s^*\|_\infty^2 \|\Delta\|_{2,\infty}^2 \|\mathbf{V}\|_{\text{F}}^2 \\
& \lesssim d^2 \delta^2 \max_{1 \leq i \leq r} \|\mathbf{u}_s^*\|_\infty^2 \|\mathbf{U}^*\|_{2,\infty}^2 \|\mathbf{V}\|_{\text{F}}^2 \\
& \lesssim \delta^2 \mu^2 r \lambda_{\max}^{*4/3} \|\mathbf{V}\|_{\text{F}}^2.
\end{aligned}$$

Regarding $p^{-1} \|\mathcal{P}_\Omega(\sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2})\|_{\mathbb{F}}^2$, we apply [YZ16, Lemma 5] (with slight modification, which we omit here for brevity) to show that: with probability exceeding $1 - O(d^{-10})$

$$\frac{1}{\sqrt{p}} \left\| \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} \right) \right\|_{\mathbb{F}} \leq \frac{3}{2} \left\| \sum_{s \in [r]} \mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} \right\|_{\mathbb{F}} = \frac{3}{2} \|\mathbf{V} \tilde{\mathbf{U}}^{*\top}\|_{\mathbb{F}} \leq \frac{3}{2} \|\tilde{\mathbf{U}}^*\| \|\mathbf{V}\|_{\mathbb{F}} \lesssim \lambda_{\max}^{*2/3} \|\mathbf{V}\|_{\mathbb{F}}.$$

under the sample size assumptin that $p \gg \mu^2 r^2 d^{-2} \log d$. Here the last inequality makes use of (90). It is self-evident that the above bounds also hold for quantities that appear in (95). Since $0 < \delta \ll 1/\sqrt{r} < 1$, we obtain

$$|\beta_1| \lesssim \delta \mu \sqrt{r} \lambda_{\max}^{*4/3} \|\mathbf{V}\|_{\mathbb{F}}^2.$$

The same upper bound holds for any other β_2 . Therefore, as long as $0 < \delta \ll 1/(\mu\sqrt{r})$ and $\kappa \asymp 1$, we have

$$|\alpha_1| \leq \frac{1}{10} \lambda_{\min}^{*4/3} \|\mathbf{V}\|_{\mathbb{F}}^2. \quad (99)$$

A.1.3 Bounding α_2

Regarding α_2 , applying [YZ16, Lemma 5] (with slight modification, which we omit here for brevity) implies that: if $p \gg \mu^2 r^2 d^{-2} \log d$, then with probability exceeding $1 - O(d^{-10})$,

$$\begin{aligned} & \left| \frac{1}{2p} \left\| \mathcal{P}_\Omega \left(\sum_{s \in [r]} (\mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} + \mathbf{u}_s^* \otimes \mathbf{v}_s \otimes \mathbf{u}_s^* + \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{v}_s) \right) \right\|_{\mathbb{F}}^2 \right. \\ & \quad \left. - \frac{1}{2} \left\| \sum_{s \in [r]} (\mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} + \mathbf{u}_s^* \otimes \mathbf{v}_s \otimes \mathbf{u}_s^* + \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{v}_s) \right\|_{\mathbb{F}}^2 \right| \\ & \leq \frac{1}{100} \left\| \sum_{s \in [r]} (\mathbf{v}_s \otimes \mathbf{u}_s^{*\otimes 2} + \mathbf{u}_s^* \otimes \mathbf{v}_s \otimes \mathbf{u}_s^* + \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{v}_s) \right\|_{\mathbb{F}}^2 = \frac{1}{100} \alpha_4 \leq \frac{1}{10} \lambda_{\min}^{*4/3} \|\mathbf{V}\|_{\mathbb{F}}^2. \end{aligned}$$

Here the last inequality arises from (92).

A.1.4 Bounding α_3

We now move on to bounding α_3 . The triangle inequality gives

$$\begin{aligned} \left| \left\langle p^{-1} \mathcal{P}_\Omega \left(\sum_{s \in [r]} \mathbf{u}_s^{\otimes 3} - \mathbf{T}^* \right), \sum_{s \in [r]} \mathbf{v}_s^{\otimes 2} \otimes \mathbf{u}_s \right\rangle \right| & \leq \sum_{s \in [r]} \left| \left\langle p^{-1} \mathcal{P}_\Omega \left(\sum_{i \in [r]} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right) \times_3 \mathbf{u}_s, \mathbf{v}_s \mathbf{v}_s^\top \right\rangle \right| \\ & \leq \max_{s \in [r]} \left\| p^{-1} \mathcal{P}_\Omega \left(\sum_{i \in [r]} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right) \times_3 \mathbf{u}_s \right\| \sum_{s \in [r]} \|\mathbf{v}_s\|_2^2 \\ & \leq \max_{s \in [r]} \left\| p^{-1} \mathcal{P}_\Omega \left(\sum_{i \in [r]} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right) \times_3 \mathbf{u}_s \right\| \|\mathbf{V}\|_{\mathbb{F}}^2. \end{aligned}$$

Recall the definitions of $\mathbf{\Delta}$ and $\mathbf{\Delta}_i$ in (94). Fix an arbitrary $s \in [r]$. From the definition of the operator norm and the triangle inequality, we can derive

$$\left\| p^{-1} \mathcal{P}_\Omega \left(\sum_{i \in [r]} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right) \times_3 \mathbf{u}_s \right\| \leq \|\mathbf{u}_s\|_2 \left\| p^{-1} \mathcal{P}_\Omega \left(\sum_{i \in [r]} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right) \right\|. \quad (100)$$

In order to upper bound $\left\| p^{-1} \mathcal{P}_\Omega \left(\sum_{i \in [r]} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right) \right\|$ as required in (100), we invoke the following simple fact, which follows immediately from the definition of the operator norm. Here and throughout, for any tensor $\mathbf{A} \in \mathbb{R}^{d \times d \times d}$ we denote

$$|\mathbf{A}| := [|A_{i,j,k}|]_{1 \leq i,j,k \leq d} \in \mathbb{R}^{d \times d \times d}.$$

Lemma A.1. *Consider any tensor $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{d \times d \times d}$ obeying $|B_{i,j,k}| \geq |A_{i,j,k}|$ for all $1 \leq i, j, k \leq d$. One has*

$$\|\mathbf{A}\| \leq \||\mathbf{A}|\| \leq \|\mathbf{B}\|. \quad (101)$$

With this lemma in mind, we are ready to derive that

$$\begin{aligned} \left\| p^{-1} \mathcal{P}_\Omega \left(\sum_{i \in [r]} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right) \right\| &\leq \left\| p^{-1} \mathcal{P}_\Omega \left(\left| \sum_{i \in [r]} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right| \right) \right\| \\ &\leq \left\| \sum_{i \in [r]} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right\|_\infty \left\| p^{-1} \mathcal{P}_\Omega(\mathbf{1}^{\otimes 3}) \right\|. \end{aligned}$$

Here, $\mathbf{1}$ stands for the all-one vector in \mathbb{R}^d . This suggests that we shall upper bound $\left\| p^{-1} \mathcal{P}_\Omega(\mathbf{1}^{\otimes 3}) \right\|$ and $\left\| \sum_{i \in [r]} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right\|_\infty$.

Given that $p \gtrsim d^{-3/2} \log^3 d$, applying Lemma D.2 indicates that with probability at least $1 - O(d^{-10})$,

$$\left\| p^{-1} \mathcal{P}_\Omega(\mathbf{1}^{\otimes 3}) - \mathbf{1}^{\otimes 3} \right\| \lesssim \frac{\log^3 d}{p} + \sqrt{\frac{d \log^5 d}{p}} \lesssim d^{3/2}. \quad (102)$$

Moreover, it is straightforward to see that $\|\mathbf{1}^{\otimes 3}\| = \|\mathbf{1}\|_2^3 = d^{3/2}$. Therefore, one has

$$\left\| p^{-1} \mathcal{P}_\Omega(\mathbf{1}^{\otimes 3}) \right\| \lesssim d^{3/2}. \quad (103)$$

Next, we turn to $\|\cdot\|_\infty$. We first expand

$$\begin{aligned} \sum_{i \in [r]} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* &= \sum_{i \in [r]} ((\Delta_i + \mathbf{u}_i^*)^{\otimes 3} - \mathbf{u}_i^{*\otimes 3}) = \sum_{i \in [r]} \Delta_i \otimes \mathbf{u}_i^{*\otimes 2} + \sum_{i \in [r]} \mathbf{u}_i^* \otimes \Delta_i \otimes \mathbf{u}_i^* + \sum_{i \in [r]} \mathbf{u}_i^{*\otimes 2} \otimes \Delta_i \\ &\quad + \sum_{i \in [r]} \Delta_i^{\otimes 2} \otimes \mathbf{u}_i^* + \sum_{i \in [r]} \Delta_i \otimes \mathbf{u}_i^* \otimes \Delta_i + \sum_{i \in [r]} \mathbf{u}_i^* \otimes \Delta_i^{\otimes 2} + \sum_{i \in [r]} \Delta_i^{\otimes 3}. \end{aligned}$$

By symmetry, it suffices to control $\left\| \sum_{i \in [r]} \mathbf{u}_i^{*\otimes 2} \otimes \Delta_i \right\|$, $\left\| \sum_{i \in [r]} \Delta_i^{\otimes 2} \otimes \mathbf{u}_i^* \right\|$ and $\left\| \sum_{i \in [r]} \Delta_i^{\otimes 3} \right\|$. Let us look at the first term. Towards this, for each $(i, j, k) \in [d]^3$, we can use the Cauchy-Schwartz inequality to control

$$\begin{aligned} \left| \left(\sum_{1 \leq s \leq r} \mathbf{u}_s^{*\otimes 2} \otimes \Delta_s \right)_{i,j,k} \right| &= \left| \sum_{1 \leq s \leq r} (\mathbf{u}_s^*)_i (\mathbf{u}_s^*)_j (\Delta_s)_k \right| \\ &\leq \left(\sum_{1 \leq s \leq r} [(\mathbf{u}_s^*)_i]^2 \right)^{1/2} \left(\sum_{1 \leq s \leq r} [(\mathbf{u}_s^*)_j]^2 [(\Delta_s)_k]^2 \right)^{1/2} \\ &\leq \|\mathbf{U}^*\|_{2,\infty} \|\Delta\|_{2,\infty} \max_{1 \leq s \leq r} \|\mathbf{u}_s^*\|_\infty \\ &\leq \delta \|\mathbf{U}^*\|_{2,\infty}^2 \max_{1 \leq s \leq r} \|\mathbf{u}_s^*\|_\infty, \end{aligned}$$

which implies that

$$\left\| \sum_{1 \leq s \leq r} \mathbf{u}_s^{*\otimes 2} \otimes \Delta_s \right\|_\infty \leq \delta \|\mathbf{U}^*\|_{2,\infty}^2 \max_{1 \leq s \leq r} \|\mathbf{u}_s^*\|_\infty \lesssim \frac{\delta \mu^{3/2} r \lambda_{\max}^*}{d^{3/2}}. \quad (104)$$

In a similar manner, we can control the remaining two terms by

$$\left\| \sum_{i \in [r]} \Delta_i^{\otimes 2} \otimes \mathbf{u}_i^* \right\| \leq \|\Delta\|_{2,\infty}^2 \max_{1 \leq s \leq r} \|\mathbf{u}_s^*\|_\infty \leq \delta^2 \|\mathbf{U}^*\|_{2,\infty}^2 \max_{1 \leq s \leq r} \|\mathbf{u}_s^*\|_\infty \leq \frac{\delta^2 \mu^{3/2} r \lambda_{\max}^*}{d^{3/2}}; \quad (105)$$

$$\left\| \sum_{i \in [r]} \Delta_i^{\otimes 3} \right\| \leq \|\Delta\|_{2,\infty}^3 \max_{1 \leq s \leq r} \|\Delta_s\|_\infty \leq \delta^3 \|\mathbf{U}^*\|_{2,\infty}^3 \leq \frac{\delta^3 \mu^{3/2} r^{3/2} \lambda_{\max}^*}{d^{3/2}}. \quad (106)$$

Recall that $0 < \delta \ll 1/r \leq 1$. Putting these together reveals that

$$\left\| \sum_{i \in [r]} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right\| \lesssim \frac{\delta \mu^{3/2} r \lambda_{\max}^*}{d^{3/2}}. \quad (107)$$

This combined with (103) yields

$$\left\| p^{-1} \mathcal{P}_\Omega \left(\sum_{i \in [r]} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right) \right\| \lesssim \delta \mu^{3/2} r \lambda_{\max}^*,$$

thus indicating that

$$\left\| p^{-1} \mathcal{P}_\Omega \left(\sum_{i \in [r]} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right) \times_3 \mathbf{u}_s \right\| \leq \|\mathbf{u}_s\|_2 \left\| p^{-1} \mathcal{P}_\Omega \left(\sum_{i \in [r]} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* \right) \right\| \lesssim \delta \mu^{3/2} r \lambda_{\max}^{4/3}, \quad (108)$$

where we use (96c) in the last step. In view of the condition that $\delta \ll 1/(\mu^{3/2} r)$ and the assumption $\kappa \asymp 1$, one has with probability greater than $1 - O(d^{-10})$,

$$|\alpha_3| \leq \frac{1}{10} \lambda_{\min}^{*4/3} \|\mathbf{V}\|_{\mathbb{F}}^2. \quad (109)$$

A.1.5 Putting all this together

Note that the above bounds hold uniformly for all \mathbf{V} . Therefore, combining upper bounds for α_i and the union bound, we conclude that with probability exceeding $1 - O(d^{-10})$,

$$\text{vec}(\mathbf{V})^\top \nabla^2 f(\mathbf{U}) \text{vec}(\mathbf{V}) \geq \alpha_4 - |\alpha_1| - |\alpha_2| - |\alpha_3| \geq \frac{1}{2} \lambda_{\min}^{*4/3} \|\mathbf{V}\|_{\mathbb{F}}^2 \quad (110)$$

$$\text{vec}(\mathbf{V})^\top \nabla^2 f(\mathbf{U}) \text{vec}(\mathbf{V}) \leq \alpha_4 + |\alpha_1| + |\alpha_2| + |\alpha_3| \leq 4 \lambda_{\max}^{*4/3} \|\mathbf{V}\|_{\mathbb{F}}^2 \quad (111)$$

as claimed.

A.2 Proof of Lemma 5.2

From (47a) and (47c), we use the triangle inequality to obtain

$$\begin{aligned} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}} &\leq \|\mathbf{U}^{t,(m)} - \mathbf{U}^t\|_{\mathbb{F}} + \|\mathbf{U}^t - \mathbf{U}^*\|_{\mathbb{F}} \\ &\leq \left(2C_1 \rho^t \mathcal{E}_{\text{local}} + 2C_2 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{\mathbb{F}} \end{aligned}$$

Similarly, we can combine (47b) and (47c) to obtain that

$$\begin{aligned} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{2,\infty} &\leq \|\mathbf{U}^{t,(m)} - \mathbf{U}^t\|_{2,\infty} + \|\mathbf{U}^t - \mathbf{U}^*\|_{\mathbb{F}} \leq \|\mathbf{U}^{t,(m)} - \mathbf{U}^t\|_{\mathbb{F}} + \|\mathbf{U}^t - \mathbf{U}^*\|_{\mathbb{F}} \\ &\leq \left(C_5 \rho^t \mathcal{E}_{\text{local}} + C_6 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty} + \left(C_3 \rho^t \mathcal{E}_{\text{local}} + C_4 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty} \\ &\leq \left((C_3 + C_5) \rho^t \mathcal{E}_{\text{local}} + (C_4 + C_6) \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty}. \end{aligned}$$

A.3 Proof of Lemma 5.3

In view of the relation (40), one has

$$\begin{aligned} \|\mathbf{U}^{t+1} - \mathbf{U}^*\|_{\mathbb{F}} &= \|\mathbf{U}^t - \eta (\nabla f_{\text{clean}}(\mathbf{U}^t) - p^{-1} \mathcal{P}_\Omega(\mathbf{E}) \times_1^{\text{seq}} \mathbf{U}^t \times_2^{\text{seq}} \mathbf{U}^t) - \mathbf{U}^*\|_{\mathbb{F}} \\ &\leq \underbrace{\|\mathbf{U}^t - \eta \nabla f_{\text{clean}}(\mathbf{U}^t) - \mathbf{U}^*\|_{\mathbb{F}}}_{=: \alpha_1} + \eta \underbrace{\|p^{-1} \mathcal{P}_\Omega(\mathbf{E}) \times_1^{\text{seq}} \mathbf{U}^t \times_2^{\text{seq}} \mathbf{U}^t\|_{\mathbb{F}}}_{=: \alpha_2}, \end{aligned}$$

which motivates us to bound α_1 and α_2 separately.

(1) We start with α_1 , towards which we find it helpful to define

$$\mathbf{U}^t(\tau) := \tau \mathbf{U}^t + (1 - \tau) \mathbf{U}^*. \quad (112)$$

Given that $\nabla f_{\text{clean}}(\mathbf{U}^*) = \mathbf{0}$ (since \mathbf{U}^* is a global optimizer of f_{clean}), we can use the fundamental theorem of calculus to obtain

$$\text{vec}(\mathbf{U}^t - \eta \nabla f_{\text{clean}}(\mathbf{U}^t) - \mathbf{U}^*) = \text{vec}(\mathbf{U}^t - \eta \nabla f_{\text{clean}}(\mathbf{U}^t) - (\mathbf{U}^* - \eta \nabla f_{\text{clean}}(\mathbf{U}^*))) \quad (113)$$

$$= \text{vec}(\mathbf{U}^t - \mathbf{U}^*) - \eta \text{vec}(\nabla f_{\text{clean}}(\mathbf{U}^t) - \nabla f_{\text{clean}}(\mathbf{U}^*)) \quad (114)$$

$$= \left(\mathbf{I}_{dr} - \eta \underbrace{\int_0^1 \nabla^2 f_{\text{clean}}(\mathbf{U}^t(\tau)) \, d\tau}_{=: \mathbf{\Gamma}} \right) \text{vec}(\mathbf{U}^t - \mathbf{U}^*). \quad (115)$$

It then follows that

$$\begin{aligned} \|\mathbf{U}^t - \eta \nabla f_{\text{clean}}(\mathbf{U}^t) - \mathbf{U}^*\|_{\text{F}}^2 &= \text{vec}(\mathbf{U}^t - \mathbf{U}^*)^\top (\mathbf{I}_{dr} - \eta \mathbf{\Gamma})^2 \text{vec}(\mathbf{U}^t - \mathbf{U}^*) \\ &\leq \|\mathbf{U}^t - \mathbf{U}^*\|_{\text{F}}^2 - 2\eta \text{vec}(\mathbf{U}^t - \mathbf{U}^*)^\top \mathbf{\Gamma} \text{vec}(\mathbf{U}^t - \mathbf{U}^*) + \eta^2 \|\mathbf{\Gamma}\|^2 \|\mathbf{U}^t - \mathbf{U}^*\|_{\text{F}}^2. \end{aligned} \quad (116)$$

From the hypothesis (47b) as well as our conditions that $\frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} + \mathcal{E}_{\text{local}} \ll \frac{1}{\mu^{3/2r}}$, we know that $\mathbf{U}^t(\tau)$ ($0 \leq \tau \leq 1$) satisfies the conditions required in Lemma 5.1. Therefore, applying Lemma 5.1 gives that

$$\begin{aligned} \text{vec}(\mathbf{U}^t - \mathbf{U}^*)^\top \mathbf{\Gamma} \text{vec}(\mathbf{U}^t - \mathbf{U}^*) &\geq \frac{1}{2} \lambda_{\min}^{*4/3} \|\mathbf{U}^t - \mathbf{U}^*\|_{\text{F}}^2, \\ \|\mathbf{\Gamma}\| &\leq 4 \lambda_{\max}^{*4/3}. \end{aligned}$$

Substitution into (116) indicates that: if $0 < \eta \leq \lambda_{\min}^{*4/3} / (32 \lambda_{\max}^{*8/3})$, then

$$\|\mathbf{U}^t - \eta \nabla f_{\text{clean}}(\mathbf{U}^t) - \mathbf{U}^*\|_{\text{F}}^2 \leq (1 - \lambda_{\min}^{*4/3} \eta + 16 \lambda_{\max}^{*8/3} \eta^2) \|\mathbf{U}^t - \mathbf{U}^*\|_{\text{F}}^2 \leq (1 - \frac{1}{2} \lambda_{\min}^{*4/3} \eta) \|\mathbf{U}^t - \mathbf{U}^*\|_{\text{F}}^2,$$

which implies that (since $1 - a/2 \geq \sqrt{1-a}$ for $0 < a < 1$)

$$\|\mathbf{U}^t - \eta \nabla f_{\text{clean}}(\mathbf{U}^t) - \mathbf{U}^*\|_{\text{F}} \leq (1 - \frac{1}{4} \lambda_{\min}^{*4/3} \eta) \|\mathbf{U}^t - \mathbf{U}^*\|_{\text{F}}. \quad (117)$$

(2) We now turn to α_2 . To simplify presentation, we shall assume that $\{E_{i,j,k}\}_{i,j,k \in [d]}$ (resp. $\{\chi_{i,j,k}\}_{i,j,k \in [d]}$) are independent random variables. Fix an arbitrary $s \in [r]$ and $m \in [d]$. The m -th entry of $\mathcal{P}_\Omega(\mathbf{E}) \times_1 \mathbf{u}_s^t \times_2 \mathbf{u}_s^t$ can be expanded as follows:

$$\begin{aligned} |(\mathcal{P}_\Omega(\mathbf{E}) \times_1 \mathbf{u}_s^t \times_2 \mathbf{u}_s^t)_m| &= \left| \mathbf{u}_s^{t\top} (\mathcal{P}_\Omega(\mathbf{E}))_{::,m} \mathbf{u}_s^t \right| \\ &\leq \underbrace{\left| \mathbf{u}_s^{t,(m)\top} (\mathcal{P}_\Omega(\mathbf{E}))_{::,m} \mathbf{u}_s^{t,(m)} \right|}_{=: \beta_1} + \underbrace{\left| (\mathbf{u}_s^t - \mathbf{u}_s^{t,(m)})^\top (\mathcal{P}_\Omega(\mathbf{E}))_{::,m} \mathbf{u}_s^t \right|}_{=: \beta_2} \\ &\quad + \underbrace{\left| \mathbf{u}_s^{t\top} (\mathcal{P}_\Omega(\mathbf{E}))_{::,m} (\mathbf{u}_s^t - \mathbf{u}_s^{t,(m)}) \right|}_{=: \beta_3} + \underbrace{\left| (\mathbf{u}_s^t - \mathbf{u}_s^{t,(m)})^\top (\mathcal{P}_\Omega(\mathbf{E}))_{::,m} (\mathbf{u}_s^t - \mathbf{u}_s^{t,(m)}) \right|}_{=: \beta_4}. \end{aligned} \quad (118)$$

Before continuing, we make the following observations: from the hypotheses (47a), (47b), (47c), as well as our assumption that $\frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} + \mathcal{E}_{\text{local}} \ll \frac{1}{\mu^{3/2r}}$, we see that the following holds for all $s \in [r]$:

$$\|\mathbf{u}_s^t - \mathbf{u}_s^*\|_2 \leq \|\mathbf{U}^t - \mathbf{U}^*\|_{\text{F}} \ll \frac{1}{r} \|\mathbf{U}^*\|_{\text{F}} \lesssim \frac{1}{\sqrt{r}} \lambda_{\max}^{*1/3}, \quad (119a)$$

$$\|\mathbf{u}_s^t - \mathbf{u}_s^*\|_\infty \leq \|\mathbf{U}^t - \mathbf{U}^*\|_{2,\infty} \ll \frac{1}{r} \|\mathbf{U}^*\|_{2,\infty} \lesssim \sqrt{\frac{\mu}{rd}} \lambda_{\max}^{*1/3}, \quad (119b)$$

$$\|\mathbf{u}_s^{t,(m)} - \mathbf{u}_s^t\|_2 \leq \|\mathbf{U}^{t,(m)} - \mathbf{U}^t\|_{\text{F}} \ll \frac{1}{r} \|\mathbf{U}^*\|_{2,\infty} \lesssim \sqrt{\frac{\mu}{rd}} \lambda_{\max}^{*1/3}, \quad (119c)$$

$$\left| (\mathbf{u}_s^{t,(m)} - \mathbf{u}_s^*)_{m'} \right| \leq \|(\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m',:}\|_2 \ll \frac{1}{r} \|\mathbf{U}^*\|_{2,\infty} \lesssim \sqrt{\frac{\mu}{rd}} \lambda_{\max}^*, \quad (119d)$$

$$\|\mathbf{u}_s^{t,(m)} - \mathbf{u}_s^*\|_2 \leq \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_F \ll \frac{1}{r} \|\mathbf{U}^*\|_F \lesssim \frac{1}{\sqrt{r}} \lambda_{\max}^{*1/3}, \quad (119e)$$

$$\|\mathbf{u}_s^t\|_2 \leq \|\mathbf{u}_s^t - \mathbf{u}_s^*\|_2 + \|\mathbf{u}_s^*\|_2 \lesssim \lambda_{\max}^{1/3}, \quad (119f)$$

$$\|\mathbf{u}_s^t\|_\infty \leq \|\mathbf{u}_s^t - \mathbf{u}_s^*\|_\infty + \|\mathbf{u}_s^*\|_\infty \lesssim \sqrt{\frac{\mu}{d}} \lambda_{\max}^{1/3}, \quad (119g)$$

$$\|\mathbf{u}_s^{t,(m)}\|_2 \leq \|\mathbf{u}_s^{t,(m)} - \mathbf{u}_s^t\|_2 + \|\mathbf{u}_s^t\|_2 \lesssim \lambda_{\max}^{*1/3}, \quad (119h)$$

$$\|\mathbf{u}_s^{t,(m)}\|_\infty \leq \|\mathbf{u}_s^{t,(m)} - \mathbf{u}_s^t\|_2 + \|\mathbf{u}_s^t\|_\infty \lesssim \sqrt{\frac{\mu}{d}} \lambda_{\max}^{1/3}, \quad (119i)$$

$$\|\mathbf{U}^t\|_F \leq \|\mathbf{U}^t - \mathbf{U}^*\|_F + \|\mathbf{U}^*\|_F \lesssim \|\mathbf{U}^*\|_F, \quad (119j)$$

$$\|\mathbf{U}^t\|_{2,\infty} \leq \|\mathbf{U}^t - \mathbf{U}^*\|_{2,\infty} + \|\mathbf{U}^*\|_{2,\infty} \lesssim \|\mathbf{U}^*\|_{2,\infty}, \quad (119k)$$

$$\|\mathbf{U}^{t,(m)}\|_F \leq \|\mathbf{U}^{t,(m)} - \mathbf{U}^t\|_F + \|\mathbf{U}^t\|_F \lesssim \|\mathbf{U}^*\|_F, \quad (119l)$$

$$\|\mathbf{U}^{t,(m)}\|_{2,\infty} \leq \|\mathbf{U}^{t,(m)} - \mathbf{U}^t\|_F + \|\mathbf{U}^t\|_{2,\infty} \lesssim \|\mathbf{U}^*\|_{2,\infty}. \quad (119m)$$

With these estimates in place, we can upper bound the above four terms in (118) separately.

- For β_1 , we note that, by construction, $\mathbf{u}^{t,(m)}$ is independent of the m -th mode-3 slice of $\mathcal{P}_\Omega(\mathbf{E})$. This tells us that

$$\mathbf{u}_s^{t,(m)\top} (\mathcal{P}_\Omega(\mathbf{E}))_{:,:,m} \mathbf{u}_s^{t,(m)} = \sum_{i,j \in [d]} (\mathbf{u}_s^{t,(m)})_i (\mathbf{u}_s^{t,(m)})_j E_{i,j,m} \chi_{i,j,m}$$

can be viewed as a sum of independent zero-mean random variables (conditional on $\mathcal{P}_{\Omega_{-m}}(\mathbf{E})$). It is straightforward to compute that

$$\begin{aligned} \max_{i,j \in [d]} \left\| (\mathbf{u}_s^{t,(m)})_i (\mathbf{u}_s^{t,(m)})_j E_{i,j,m} \chi_{i,j,m} \right\|_{\psi_1} &\lesssim \sigma \|\mathbf{u}_s^{t,(m)}\|_\infty^2 =: L, \\ \sum_{i,j \in [d]} \mathbb{E} \left[(\mathbf{u}_s^{t,(m)})_i^2 (\mathbf{u}_s^{t,(m)})_j^2 E_{i,j,m}^2 \chi_{i,j,m}^2 \right] &\lesssim \sigma^2 p \|\mathbf{u}_s^{t,(m)}\|_2^4 \lesssim \sigma^2 \lambda_{\max}^{*2/3} p \|\mathbf{u}_s^{t,(m)}\|_2^2 =: V, \end{aligned}$$

where $\|\cdot\|_{\psi_1}$ denotes the sub-exponential norm and we use (119) in the last inequality. Applying the matrix Bernstein inequality [Kol11, Corollary 2.1] yields that: with probability $1 - O(d^{-20})$,

$$\begin{aligned} \left| \mathbf{u}_s^{t,(m)\top} (\mathcal{P}_\Omega(\mathbf{E}))_{:,:,m} \mathbf{u}_s^{t,(m)} \right| &\lesssim L \log^2 d + \sqrt{V \log d} \\ &\lesssim \sigma \log^2 d \|\mathbf{u}_s^{t,(m)}\|_\infty^2 + \sigma \lambda_{\max}^{*1/3} \sqrt{p \log d} \|\mathbf{u}_s^{t,(m)}\|_2, \end{aligned} \quad (120)$$

- For β_2 , we first invoke [CW15, lemma 11] to demonstrate that: with probability at least $1 - O(d^{-11})$,

$$\max_{m \in [d]} \left\| (\mathcal{P}_\Omega(\mathbf{E}))_{:,:,m} \right\| \lesssim \sigma (\sqrt{dp} + \log d). \quad (121)$$

Next, it is seen that

$$\begin{aligned} \left| (\mathbf{u}_s^t - \mathbf{u}_s^{t,(m)})^\top (\mathcal{P}_\Omega(\mathbf{E}))_{:,:,m} \mathbf{u}_s^t \right| &\leq \left\| (\mathcal{P}_\Omega(\mathbf{E}))_{:,:,m} \right\| \|\mathbf{u}_s^t - \mathbf{u}_s^{t,(m)}\|_2 \|\mathbf{u}_s^t\|_2 \\ &\lesssim \sigma (\sqrt{dp} + \log d) \|\mathbf{U}^t - \mathbf{U}^{t,(m)}\|_F \|\mathbf{u}_s^t\|_2 \\ &\lesssim \sigma (\sqrt{dp} + \log d) \left(C_5 \rho^t \mathcal{E}_{\text{local}} + C_6 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty} \|\mathbf{u}_s^t\|_2 \\ &\ll \frac{\lambda_{\min}^{*4/3}}{\sqrt{d}} \left(C_5 \rho^t \mathcal{E}_{\text{local}} + C_6 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{u}_s^t\|_2, \end{aligned}$$

where the last inequality follows from the noise condition. Clearly, the above bound holds for β_3 as well.

- Regarding β_4 , it is easily seen that $\beta_4 \ll \beta_2$, given that $\|\mathbf{U}^t - \mathbf{U}^{t,(m)}\|_F \ll \|\mathbf{U}^*\|_{2,\infty}/r \leq \lambda_{\min}^*$ holds according to (119) and $\kappa \asymp 1$.
- Taking together the above bounds and substituting them into (118), we obtain

$$\begin{aligned} |(\mathcal{P}_\Omega(\mathbf{E}) \times_1 \mathbf{u}_s^t \times_2 \mathbf{u}_s^t)_m| &\lesssim \sigma \log^2 d \|\mathbf{u}_s^{t,(m)}\|_\infty^2 + \sigma \lambda_{\max}^{*1/3} \sqrt{p \log d} \|\mathbf{u}_s^{t,(m)}\|_2 \\ &\quad + o(1) \frac{\lambda_{\min}^{*4/3}}{\sqrt{d}} \left(C_5 \rho^t \mathcal{E}_{\text{local}} + C_6 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{u}_s^t\|_2. \end{aligned}$$

Recognizing that this holds for any $m \in [d]$ and $s \in [r]$, one can sum over m and s to deduce that

$$\begin{aligned} \|\mathcal{P}_\Omega(\mathbf{E}) \times_1^{\text{seq}} \mathbf{U}^t \times_2^{\text{seq}} \mathbf{U}^t\|_F &\stackrel{(i)}{\lesssim} \frac{\sigma \lambda_{\max}^{*2/3} \mu \sqrt{r} \log^2 d}{\sqrt{d}} + \sigma \lambda_{\max}^{*1/3} \sqrt{dp \log d} \|\mathbf{U}^{t,(m)}\|_F \\ &\quad + o(1) \lambda_{\min}^{*4/3} \left(C_5 \rho^t \mathcal{E}_{\text{local}} + C_6 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^t\|_F \\ &\leq C \sigma \lambda_{\min}^{*1/3} \sqrt{dp \log d} \|\mathbf{U}^*\|_F \\ &\quad + o(1) \lambda_{\min}^{*4/3} \left(C_5 \rho^t \mathcal{E}_{\text{local}} + C_6 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_F, \end{aligned} \quad (122)$$

for some absolute constant $C > 0$, where (i) is true due to (119); the last inequality follows from the fact that $\|\mathbf{U}^*\|_F \geq \lambda_{\min}^{*1/3} \sqrt{r}$ as well as the assumptions that $p \gg \mu d^{-2} \log^3 d$ and $\kappa \asymp 1$.

(3) Combining (117) and (122) yields that: with probability at least $1 - O(d^{-10})$,

$$\begin{aligned} \|\mathbf{U}^{t+1} - \mathbf{U}^*\|_F &\leq \left(1 - \frac{1}{4} \lambda_{\min}^{*4/3} \eta\right) \|\mathbf{U}^t - \mathbf{U}^*\|_F + C \eta \sigma \lambda_{\min}^{*1/3} \sqrt{\frac{d \log d}{p}} \|\mathbf{U}^*\|_F \\ &\quad + o(1) \lambda_{\min}^{*4/3} \eta \left(C_5 \rho^t \mathcal{E}_{\text{local}} + C_6 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_F \\ &\leq \left(1 - \frac{1}{4} \lambda_{\min}^{*4/3} \eta\right) \left(C_1 \rho^t \mathcal{E}_{\text{local}} + C_2 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_F + C \eta \sigma \lambda_{\min}^{*1/3} \sqrt{\frac{d \log d}{p}} \|\mathbf{U}^*\|_F \\ &\quad + o(1) \lambda_{\min}^{*4/3} \eta \left(C_5 \rho^t \mathcal{E}_{\text{local}} + C_6 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_F \\ &\leq \left(1 - \frac{1}{5} \lambda_{\min}^{*4/3} \eta\right) C_1 \rho^t \mathcal{E}_{\text{local}} \|\mathbf{U}^*\|_F + \left(\left(1 - \frac{1}{5} \lambda_{\min}^{*4/3} \eta\right) C_2 + C \eta \lambda_{\min}^{*3/4} \right) \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \|\mathbf{U}^*\|_F \\ &\leq C_1 \rho^{t+1} \mathcal{E}_{\text{local}} \|\mathbf{U}^*\|_F + C_2 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \|\mathbf{U}^*\|_F, \end{aligned}$$

with the proviso that $0 < \eta \leq \lambda_{\min}^{*4/3}/(32\lambda_{\max}^{*8/3})$, $1 - (\lambda_{\min}^{*4/3}/5)\eta \leq \rho < 1$ and that C_2 is sufficiently large.

A.4 Proof of Lemma 5.4

Fix an arbitrary $m \in [r]$. From the definition of $f^{(m)}$ in (45) and (46), we can show that

$$\begin{aligned} \mathbf{U}^{t+1} - \mathbf{U}^{t+1,(m)} &= \mathbf{U}^t - \eta \nabla f(\mathbf{U}^t) - \left(\mathbf{U}^{t,(m)} - \eta \nabla f^{(m)}(\mathbf{U}^{t,(m)}) \right) \\ &= \mathbf{U}^t - \mathbf{U}^{t,(m)} - \eta \left(\nabla f(\mathbf{U}^t) - \nabla f(\mathbf{U}^{t,(m)}) \right) - \eta \left(\nabla f(\mathbf{U}^{t,(m)}) - \nabla f^{(m)}(\mathbf{U}^{t,(m)}) \right) \end{aligned}$$

$$\begin{aligned}
&= \underbrace{\mathbf{U}^t - \mathbf{U}^{t,(m)} - \eta \left(\nabla f_{\text{clean}}(\mathbf{U}^t) - \nabla f_{\text{clean}}(\mathbf{U}^{t,(m)}) \right)}_{=: \alpha_1} \\
&\quad - \underbrace{\eta \left((p^{-1}\mathcal{P}_{\Omega_m} - \mathcal{P}_m) \left(\sum_{s \in [r]} (\mathbf{u}_s^{t,(m)})^{\otimes 3} - \mathbf{T}^* \right) \right)}_{=: \alpha_2} \times_1^{\text{seq}} \mathbf{U}^{t,(m)} \times_2^{\text{seq}} \mathbf{U}^{t,(m)} \\
&\quad + \underbrace{\eta p^{-1} \mathcal{P}_{\Omega_m}(\mathbf{E}) \times_1^{\text{seq}} \mathbf{U}^{t,(m)} \times_2^{\text{seq}} \mathbf{U}^{t,(m)}}_{=: \alpha_3} + \underbrace{\eta p^{-1} \{ \mathcal{P}_{\Omega}(\mathbf{E}) \times_1^{\text{seq}} \mathbf{U}^t \times_2^{\text{seq}} \mathbf{U}^t - \mathcal{P}_{\Omega}(\mathbf{E}) \times_1^{\text{seq}} \mathbf{U}^{t,(m)} \times_2^{\text{seq}} \mathbf{U}^{t,(m)} \}}_{=: \alpha_4}.
\end{aligned} \tag{123}$$

Before proceeding to bounding these terms, we pause to define

$$\Delta_{\mathbf{T}}^{t,(m)} := \sum_{s \in [r]} (\mathbf{u}_s^{t,(m)})^{\otimes 3} - \mathbf{T}^*. \tag{124}$$

From (119) and hypothesis (47d), one can applying a similar argument as in (107) to find that

$$\|\Delta_{\mathbf{T}}^{t,(m)}\|_{\infty} \lesssim \max_{s \in [r]} \|\mathbf{u}_s^*\|_{\infty} \|\mathbf{U}^*\|_{2,\infty} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{2,\infty} \lesssim \frac{\mu \sqrt{r} \lambda_{\max}^{*2/3}}{d} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{2,\infty}. \tag{125}$$

Now we begin to bound the terms in (123) separately.

(1) We start with α_1 . For any $0 \leq \tau \leq 1$, define

$$\mathbf{U}^{t,(m)}(\tau) := \tau \mathbf{U}^t + (1 - \tau) \mathbf{U}^{t,(m)}.$$

The fundamental theorem of calculus yields

$$\mathbf{U}^t - \mathbf{U}^{t,(m)} - \eta \left(\nabla f_{\text{clean}}(\mathbf{U}^t) - \nabla f_{\text{clean}}(\mathbf{U}^{t,(m)}) \right) = \underbrace{\left(\mathbf{I}_{dr} - \eta \int_0^1 \nabla^2 f_{\text{clean}}(\mathbf{U}^{t,(m)}(\tau)) d\tau \right)}_{=: \Gamma} (\mathbf{U}^t - \mathbf{U}^{t,(m)}).$$

By Lemma 5.2 and our assumptions on the noise, we know that $\mathbf{U}^{t,(m)}(\tau)$ ($0 \leq \tau \leq 1$) satisfies the conditions in Lemma 5.1 for any $\tau \in [0, 1]$. Applying the same argument as the one used to bound $\|\mathbf{U}^t - \eta \nabla f_{\text{clean}}(\mathbf{U}^t) - \mathbf{U}^*\|_{\mathbb{F}}$ in Lemma 5.3, we show that

$$\left\| \mathbf{U}^t - \mathbf{U}^{t,(m)} - \eta \left(\nabla f(\mathbf{U}^t) - \nabla f(\mathbf{U}^{t,(m)}) \right) \right\|_2 \leq \left(1 - \frac{1}{4} \lambda_{\min}^{*4/3} \eta \right) \|\mathbf{U}^t - \mathbf{U}^{t,(m)}\|_{\mathbb{F}}, \tag{126}$$

provided that $0 < \eta \leq \lambda_{\min}^{*4/3} / (32 \lambda_{\max}^{*8/3})$.

In what follows, we shall assume that $\{E_{i,j,k}\}_{i,j,k \in [d]}$ (resp. $\{\chi_{i,j,k}\}_{i,j,k \in [d]}$) are independent random variables in order to simplify presentation.

(2) The next step is to bound α_2 . For notational simplicity, define

$$\begin{aligned}
\mathbf{V}^{t,(m)} &:= \left((p^{-1}\mathcal{P}_{\Omega_m} - \mathcal{P}_m) (\Delta_{\mathbf{T}}^{t,(m)}) \right) \times_1^{\text{seq}} \mathbf{U}^{t,(m)} \times_2^{\text{seq}} \mathbf{U}^{t,(m)}; \\
\mathbf{v}_s^{t,(m)} &:= \left((p^{-1}\mathcal{P}_{\Omega_m} - \mathcal{P}_m) (\Delta_{\mathbf{T}}^{t,(m)}) \right) \times_1 \mathbf{u}_s^{t,(m)} \times_2 \mathbf{u}_s^{t,(m)}, \quad s \in [r].
\end{aligned}$$

In order to control the Frobenius norm of $\mathbf{V}^{t,(m)}$, we shall start by considering the m -th row of $\mathbf{V}^{t,(m)}$. In view of the definitions of \mathcal{P}_{Ω_m} and \mathcal{P}_m (cf. Appendix 5.1.3), we can expand

$$\mathbf{V}_{m,:}^{t,(m)} = \sum_{i,j \in [d]} \sum_{s \in [r]} (p^{-1} \chi_{i,j,m} - 1) (\Delta_{\mathbf{T}}^{t,(m)})_{i,j,m} (\mathbf{u}_s^{t,(m)})_i (\mathbf{u}_s^{t,(m)})_j \mathbf{e}_s^{\top}.$$

We recognize that $\{\mathbf{u}_s^{t,(m)}\}_{s=1}^r$ is independent of Ω_m , making it convenient for us to upper bound $\|\mathbf{V}_{m,:}^{t,(m)}\|_2$. Specifically, for any $i, j \in [d]$, from (119) and (125) we have

$$\begin{aligned} \left\| \sum_{s \in [r]} (\Delta_{\mathbf{T}}^{t,(m)})_{i,j,m} (\mathbf{u}_s^{t,(m)})_i (\mathbf{u}_s^{t,(m)})_j (p^{-1} \chi_{i,j,m} - 1) \mathbf{e}_s^\top \right\|_2 &\leq \frac{1}{p} \|\Delta_{\mathbf{T}}^{t,(m)}\|_\infty \max_{s \in [r]} \|\mathbf{u}_s^{t,(m)}\|_\infty \|\mathbf{U}^{t,(m)}\|_{2,\infty} \\ &\lesssim \frac{\mu^2 r \lambda_{\max}^{*4/3}}{d^2 p} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{2,\infty} =: L_1. \end{aligned}$$

In addition, it is easy to verify that $\mathbb{E}[\mathbf{V}_{m,:}^{t,(m)}] = \mathbf{0}$ and

$$\begin{aligned} &\sum_{i,j \in [d]} \mathbb{E} \left[\left\| \sum_{s \in [r]} (\Delta_{\mathbf{T}}^{t,(m)})_{i,j,m} (\mathbf{u}_s^{t,(m)})_i (\mathbf{u}_s^{t,(m)})_j (p^{-1} \chi_{i,j,m} - 1) \mathbf{e}_s^\top \right\|_2^2 \right] \\ &= \sum_{i,j \in [d]} \sum_{s \in [r]} (\Delta_{\mathbf{T}}^{t,(m)})_{i,j,m}^2 (\mathbf{u}_s^{t,(m)})_i^2 (\mathbf{u}_s^{t,(m)})_j^2 \mathbb{E}[(p^{-1} \chi_{i,j,m} - 1)^2] \\ &\leq \frac{1}{p} \|\Delta_{\mathbf{T}}^{t,(m)}\|_\infty^2 \max_{s \in [r]} \|\mathbf{u}_s^{t,(m)}\|_2^2 \|\mathbf{U}^{t,(m)}\|_F^2 \lesssim \frac{\mu^2 r^2 \lambda_{\max}^{*8/3}}{d^2 p} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{2,\infty}^2 =: V_1, \end{aligned}$$

where the last inequality holds due to (119) and (125). We then apply the matrix Bernstein inequality to yield that: with probability exceeding $1 - O(d^{-20})$,

$$\begin{aligned} \|\mathbf{V}_{m,:}^{t,(m)}\|_2 &\lesssim L_1 \log d + \sqrt{V_1 \log d} \lesssim \left\{ \frac{\mu^2 r \lambda_{\max}^{*4/3} \log d}{d^2 p} + \frac{\mu r \lambda_{\max}^{*4/3} \sqrt{\log d}}{d \sqrt{p}} \right\} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_\infty \\ &\asymp \frac{\mu r \lambda_{\max}^{*4/3} \sqrt{\log d}}{d \sqrt{p}} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_\infty, \end{aligned} \quad (127)$$

where the last step holds as long as $p \gg \mu^2 d^{-2} \log d$.

Next, we turn to the k -th row of $\mathbf{V}^{t,(m)}$ for any $k \neq m$. For each $s \in [r]$, we have

$$\begin{aligned} (\mathbf{v}_s^{t,(m)})_k &= (\mathbf{u}_s^{t,(m)})_m \sum_{j \in [d]} (\Delta_{\mathbf{T}}^{t,(m)})_{m,j,k} (\mathbf{u}_s^{t,(m)})_j (p^{-1} \chi_{m,j,k} - 1) \\ &\quad + (\mathbf{u}_s^{t,(m)})_m \sum_{i:i \neq m} (\Delta_{\mathbf{T}}^{t,(m)})_{i,m,k} (\mathbf{u}_s^{t,(m)})_i (p^{-1} \chi_{i,m,k} - 1). \end{aligned}$$

Similar to the proof of Lemma D.9, we can show that with probability at least $1 - O(d^{-20})$,

$$\begin{aligned} &\sum_{s \in [r]} \sum_{k:k \neq m} \left((\mathbf{u}_s^{t,(m)})_m \sum_{j \in [d]} (\Delta_{\mathbf{T}}^{t,(m)})_{m,j,k} (\mathbf{u}_s^{t,(m)})_j (p^{-1} \chi_{m,j,k} - 1) \right)^2 \\ &\lesssim \sum_{s \in [r]} \sum_{k:k \neq m} \mathbb{E} \left[\left((\mathbf{u}_s^{t,(m)})_m \sum_{j \in [d]} (\Delta_{\mathbf{T}}^{t,(m)})_{m,j,k} (\mathbf{u}_s^{t,(m)})_j (p^{-1} \chi_{m,j,k} - 1) \right)^2 \right] \\ &\lesssim \frac{1}{p} \sum_{s \in [r]} (\mathbf{u}_s^{t,(m)})_m^2 \sum_{j,k \in [d]} (\Delta_{\mathbf{T}}^{t,(m)})_{m,j,k}^2 (\mathbf{u}_s^{t,(m)})_j^2 \\ &\leq \frac{d}{p} \max_{s \in [r]} \|\mathbf{u}_s^{t,(m)}\|_\infty^2 \|\Delta_{\mathbf{T}}^{t,(m)}\|_\infty^2 \|\mathbf{U}^{t,(m)}\|_F^2 \lesssim \frac{\mu^3 r^2 \lambda_{\max}^{*8/3}}{d^2 p} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{2,\infty}^2, \end{aligned}$$

where the last inequality follows from (119), and (125). It is easily seen that the bound also holds for the summation over $i \neq m$. We can then use Cauchy-Schwartz to arrive at

$$\sum_{k:k \neq m} \|\mathbf{V}_{k,:}^{t,(m)}\|_2^2 = \sum_{s \in [r]} \sum_{k:k \neq m} (\mathbf{v}_s^{t,(m)})_k^2 \lesssim \frac{\mu^3 r^2 \lambda_{\max}^{*8/3}}{d^2 p} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{2,\infty}^2. \quad (128)$$

Combining (128) with (127) and invoking the union bound, we conclude that: with probability exceeding $1 - O(d^{-20})$,

$$\left\| \left((p^{-1} \mathcal{P}_{\Omega_m} - \mathcal{P}_m) (\Delta_{\mathbf{T}}^{t,(m)}) \right) \times_1^{\text{seq}} \mathbf{U}^{t,(m)} \times_2^{\text{seq}} \mathbf{U}^{t,(m)} \right\|_{\text{F}} = \|\mathbf{V}^{t,(m)}\|_{\text{F}} \leq C \frac{\mu^{3/2} r \lambda_{\min}^{*4/3} \sqrt{\log d}}{d \sqrt{p}} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{2,\infty}. \quad (129)$$

for some absolute constant $C > 0$, where we use the assumption that $\kappa \asymp 1$.

(3) For α_3 , following a similar argument for α_2 , we define

$$\begin{aligned} \mathbf{W}^{t,(m)} &:= \mathcal{P}_{\Omega_m}(\mathbf{E}) \times_1^{\text{seq}} \mathbf{U}_s^{t,(m)} \times_2^{\text{seq}} \mathbf{U}_s^{t,(m)} \\ \mathbf{w}_s^{t,(m)} &:= \mathcal{P}_{\Omega_m}(\mathbf{E}) \times_1 \mathbf{u}_s^{t,(m)} \times_2 \mathbf{u}_s^{t,(m)} \end{aligned}$$

for each $s \in [r]$. The m -th row of $\mathbf{W}^{t,(m)}$ is a sum of independent zero-mean random vectors:

$$\mathbf{W}_{m,:}^{t,(m)} = \sum_{i,j \in [d]} \sum_{s \in [r]} (\mathbf{u}_s^{t,(m)})_i (\mathbf{u}_s^{t,(m)})_j E_{i,j,m} \chi_{i,j,m} \mathbf{e}_s^\top.$$

With (119) in place, it is easy to verify that

$$\begin{aligned} \max_{i,j \in [d]} \left\| \sum_{s \in [r]} (\mathbf{u}_s^{t,(m)})_i (\mathbf{u}_s^{t,(m)})_j E_{i,j,m} \chi_{i,j,m} \mathbf{e}_s^\top \right\|_{\psi_1} &\leq \sigma \max_{s \in [r]} \|\mathbf{u}_s^{t,(m)}\|_{\infty} \|\mathbf{U}^{t,(m)}\|_{2,\infty} \\ &\lesssim \sigma \sqrt{\frac{\mu}{d}} \lambda_{\max}^{*1/3} \|\mathbf{U}^*\|_{2,\infty} =: L_2 \end{aligned}$$

and

$$\begin{aligned} \sum_{i,j \in [d]} \mathbb{E} \left[\left\| \sum_{s \in [r]} (\mathbf{u}_s^{t,(m)})_i (\mathbf{u}_s^{t,(m)})_j E_{i,j,m} \chi_{i,j,m} \mathbf{e}_s^\top \right\|_2^2 \right] &= \sum_{i,j \in [d]} \sum_{s \in [r]} (\mathbf{u}_s^{t,(m)})_i^2 (\mathbf{u}_s^{t,(m)})_j^2 \mathbb{E} [(E_{i,j,m} \chi_{i,j,m})^2] \\ &\leq \sigma^2 p \max_{s \in [r]} \|\mathbf{u}_s^{t,(m)}\|_2^2 \|\mathbf{U}^{t,(m)}\|_{\text{F}}^2 \\ &\lesssim \sigma^2 d p \lambda_{\max}^{*2/3} \|\mathbf{U}^*\|_{2,\infty}^2 =: V_2, \end{aligned}$$

where we have used (119) and the fact that $\|\mathbf{U}^*\|_{\text{F}} \leq \sqrt{d} \|\mathbf{U}^*\|_{2,\infty}$. Apply the matrix Bernstein inequality to reveal that: with probability at least $1 - O(d^{-20})$,

$$\begin{aligned} \left\| \mathbf{W}_{m,:}^{t,(m)} \right\|_2 &= \left\| \sum_{i,j \in [d]} \sum_{s \in [r]} (\mathbf{u}_s^{t,(m)})_i (\mathbf{u}_s^{t,(m)})_j E_{i,j,m} \chi_{i,j,m} \mathbf{e}_s^\top \right\|_2 \\ &\lesssim L_2 \log^2 d + \sqrt{V_2 \log d} \asymp \sigma \lambda_{\max}^{*1/3} \sqrt{d p \log d} \|\mathbf{U}^*\|_{2,\infty}, \end{aligned} \quad (130)$$

where the last inequality holds as long as $p \gg \mu d^{-2} \log^3 d$.

As for the other rows, we have

$$(\mathbf{w}_s^{t,(m)})_k = (\mathbf{u}_s^{t,(m)})_m \sum_{j \in [d]} (\mathbf{u}_s^{t,(m)})_j E_{m,j,k} \chi_{m,j,k} + (\mathbf{u}_s^{t,(m)})_m \sum_{i:i \neq m} (\mathbf{u}_s^{t,(m)})_i E_{i,m,k} \chi_{i,m,k}$$

for each $s \in [r]$. Arguing similarly as in the proof of Lemma D.10, we have with probability at least $1 - O(d^{-20})$,

$$\begin{aligned} \sum_{s \in [r]} \sum_{k:k \neq m} \left((\mathbf{u}_s^{t,(m)})_m \sum_{j \in [d]} (\mathbf{u}_s^{t,(m)})_j E_{m,j,k} \chi_{m,j,k} \right)^2 \\ \lesssim \sum_{s \in [r]} \sum_{k:k \neq m} \mathbb{E} \left[\left((\mathbf{u}_s^{t,(m)})_m \sum_{j \in [d]} (\mathbf{u}_s^{t,(m)})_j (E_{m,j,k} \chi_{m,j,k}) \right)^2 \right] \end{aligned}$$

$$\begin{aligned}
&\leq \sigma^2 dp \sum_{s \in [r]} (\mathbf{u}_s^{t,(m)})_m^2 \sum_{j \in [d]} (\mathbf{u}_s^{t,(m)})_j^2 \\
&\leq \sigma^2 dp \max_{s \in [r]} \|\mathbf{u}_s^{t,(m)}\|_2^2 \|\mathbf{U}^{t,(m)}\|_{2,\infty}^2 \lesssim \sigma^2 \lambda_{\max}^{*2/3} dp \|\mathbf{U}^*\|_{2,\infty}^2,
\end{aligned}$$

where the last inequality follows from (119). Additionally, the summation over $\{i : i \neq m\}$ can be controlled using the same argument. Therefore, we use Cauchy-Schwarz to find that

$$\sum_{k:k \neq m} \|\mathbf{W}_{k,:}^{t,(m)}\|_2^2 \lesssim \sigma^2 \lambda_{\max}^{*2/3} dp \|\mathbf{U}^*\|_{2,\infty}^2. \quad (131)$$

Combined with (130) and the assumption that $\kappa \asymp 1$, we obtain that

$$\left\| p^{-1} \mathcal{P}_{\Omega_m}(\mathbf{E}) \times_1^{\text{seq}} \mathbf{U}^{t,(m)} \times_2^{\text{seq}} \mathbf{U}^{t,(m)} \right\|_{\mathbb{F}} = \|\mathbf{W}^{t,(m)}\|_{\mathbb{F}} \leq \tilde{C} \sigma \lambda_{\min}^{*1/3} \sqrt{dp \log d} \|\mathbf{U}^*\|_{2,\infty}, \quad (132)$$

for some absolute constant $\tilde{C} > 0$.

(4) Regarding α_4 , we use the triangle inequality to show that for each $s \in [r]$,

$$\begin{aligned}
&\left\| \mathcal{P}_{\Omega}(\mathbf{E}) \times_1 \mathbf{u}_s^t \times_2 \mathbf{u}_s^t - \mathcal{P}_{\Omega}(\mathbf{E}) \times_1 \mathbf{u}_s^{t,(m)} \times_2 \mathbf{u}_s^{t,(m)} \right\|_2 \\
&\leq \left\| \mathcal{P}_{\Omega}(\mathbf{E}) \times_1 (\mathbf{u}_s^t - \mathbf{u}_s^{t,(m)}) \times_2 \mathbf{u}_s^{t,(m)} \right\|_2 + \left\| \mathcal{P}_{\Omega}(\mathbf{E}) \times_1 \mathbf{u}_s^{t,(m)} \times_2 (\mathbf{u}_s^t - \mathbf{u}_s^{t,(m)}) \right\|_2 \\
&\quad + \left\| \mathcal{P}_{\Omega}(\mathbf{E}) \times_1 (\mathbf{u}_s^t - \mathbf{u}_s^{t,(m)}) \times_2 (\mathbf{u}_s^t - \mathbf{u}_s^{t,(m)}) \right\|_2 \\
&\leq 2 \|\mathcal{P}_{\Omega}(\mathbf{E})\| \|\mathbf{u}_s^t - \mathbf{u}_s^{t,(m)}\|_2 \|\mathbf{u}_s^{t,(m)}\|_2 + \|\mathcal{P}_{\Omega}(\mathbf{E})\| \|\mathbf{u}_s^t - \mathbf{u}_s^{t,(m)}\|_2^2 \\
&\lesssim \|\mathcal{P}_{\Omega}(\mathbf{E})\| \|\mathbf{u}_s^t - \mathbf{u}_s^{t,(m)}\|_2 \|\mathbf{u}_s^{t,(m)}\|_2,
\end{aligned}$$

where the last line follows from (119). From Corollary D.3, we can further upper bound

$$\left\| \mathcal{P}_{\Omega}(\mathbf{E}) \times_1 \mathbf{u}_s^t \times_2 \mathbf{u}_s^t - \mathcal{P}_{\Omega}(\mathbf{E}) \times_1 \mathbf{u}_s^{t,(m)} \times_2 \mathbf{u}_s^{t,(m)} \right\|_2 \leq \sigma (\sqrt{dp} + \log d) \log^{5/2} d \|\mathbf{u}_s^t - \mathbf{u}_s^{t,(m)}\|_2 \|\mathbf{u}_s^{t,(m)}\|_2.$$

As a result, we sum over $s \in [r]$ and use the Cauchy-Schwartz inequality to derive

$$\begin{aligned}
&\frac{1}{p} \left\| \mathcal{P}_{\Omega}(\mathbf{E}) \times_1^{\text{seq}} \mathbf{U}^t \times_2^{\text{seq}} \mathbf{U}^t - \mathcal{P}_{\Omega}(\mathbf{E}) \times_1^{\text{seq}} \mathbf{U}^{t,(m)} \times_2^{\text{seq}} \mathbf{U}^{t,(m)} \right\|_{\mathbb{F}} \\
&\lesssim \frac{\sigma}{p} (\sqrt{dp} + \log d) \log^{5/2} d \|\mathbf{U}^t - \mathbf{U}^{t,(m)}\|_{\mathbb{F}} \|\mathbf{U}^{t,(m)}\|_{\mathbb{F}} \\
&\lesssim \frac{\sigma}{p} (\sqrt{dp} + \log d) \log^{5/2} d \left(C_5 \rho^t \mathcal{E}_{\text{local}} + C_6 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty}^2 \\
&\ll \lambda_{\min}^{*4/3} \left(C_5 \rho^t \mathcal{E}_{\text{local}} + C_6 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty}, \quad (133)
\end{aligned}$$

where the last step arises from conditions that $\frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \ll \frac{1}{\mu^{3/2r}}$ and $\kappa \asymp 1$.

(5) Taking (126), (129), (132) and (133) together, we can invoke the sample size assumption that $p \gg \mu^3 r^2 d^{-2} \log^3 d$ and the union bound to show that: with probability greater than $1 - O(d^{-10})$ one has

$$\begin{aligned}
&\|\mathbf{U}^{t+1} - \mathbf{U}^{t+1,(m)}\|_{\mathbb{F}} \\
&\leq \left(1 - \frac{1}{4} \lambda_{\min}^{*4/3} \eta\right) \|\mathbf{U}^t - \mathbf{U}^{t,(m)}\|_{\mathbb{F}} + C \frac{\mu^{3/2r} \lambda_{\min}^{*4/3} \sqrt{\log d}}{d \sqrt{p}} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{2,\infty} \\
&\quad + \tilde{C} \eta \sigma \lambda_{\min}^{*1/3} \sqrt{\frac{\mu r d \log d}{p}} \|\mathbf{U}^*\|_{2,\infty} + o(1) \lambda_{\min}^{*4/3} \eta \left(C_5 \rho^t \mathcal{E}_{\text{local}} + C_6 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty}
\end{aligned}$$

$$\begin{aligned}
&\leq \left(1 - \frac{1}{4}\lambda_{\min}^{*4/3}\eta\right) \left(C_5\rho^t\mathcal{E}_{\text{local}} + C_6\frac{\sigma}{\lambda_{\min}^*}\sqrt{\frac{d\log d}{p}}\right) \|\mathbf{U}^*\|_{2,\infty} + \tilde{C}\eta\sigma\lambda_{\min}^{*1/3}\sqrt{\frac{d\log d}{p}}\|\mathbf{U}^*\|_{2,\infty} \\
&\quad + C\eta\lambda_{\max}^{*4/3}\frac{\mu^{3/2}r\sqrt{\log d}}{d\sqrt{p}} \left(2C_3\rho^t\mathcal{E}_{\text{local}} + 2C_4\frac{\sigma}{\lambda_{\min}^*}\sqrt{\frac{d\log d}{p}}\right) \|\mathbf{U}^*\|_{2,\infty} \\
&\quad + o(1)\lambda_{\min}^{*4/3}\eta \left(C_5\rho^t\mathcal{E}_{\text{local}} + C_6\frac{\sigma}{\lambda_{\min}^*}\sqrt{\frac{\mu rd\log d}{p}}\right) \|\mathbf{U}^*\|_{2,\infty} \\
&\leq \left(1 - \frac{1}{5}\lambda_{\min}^{*4/3}\eta\right)C_5\rho^t\mathcal{E}_{\text{local}}\|\mathbf{U}^*\|_{2,\infty} + \left(\left(1 - \frac{1}{5}\lambda_{\min}^{*4/3}\eta\right)C_6 + \tilde{C}\lambda_{\min}^{*4/3}\eta\right)\frac{\sigma}{\lambda_{\min}^*}\sqrt{\frac{d\log d}{p}}\|\mathbf{U}^*\|_{2,\infty} \\
&\leq \left(1 - \frac{1}{5}\lambda_{\min}^{*4/3}\eta\right)C_5\rho^t\mathcal{E}_{\text{local}}\|\mathbf{U}^*\|_{2,\infty} + C_6\frac{\sigma}{\lambda_{\min}^*}\sqrt{\frac{d\log d}{p}}\|\mathbf{U}^*\|_{2,\infty} \\
&\leq \left(C_5\rho^{t+1}\mathcal{E}_{\text{local}} + C_6\frac{\sigma}{\lambda_{\min}^*}\sqrt{\frac{d\log d}{p}}\right)\|\mathbf{U}^*\|_{2,\infty},
\end{aligned}$$

provided $0 < \eta \leq \lambda_{\min}^{*4/3}/(32\lambda_{\max}^{*8/3})$, $1 - (\lambda_{\min}^{*4/3}/5)\eta \leq \rho < 1$ and C_6 is sufficiently large.

A.5 Proof of Lemma 5.5

Fix an arbitrary $m \in [d]$. Recall our notation of $\Delta_{\mathbf{T}}^{t,(m)}$ in (124). To simplify presentation, we further define

$$\widehat{\mathbf{U}}^{t+1,(m)} := \mathbf{U}^{t,(m)} - \eta \left(p^{-1}\mathcal{P}_{\Omega_{-m}}(\Delta_{\mathbf{T}}^{t,(m)} - \mathbf{E}) \times_1^{\text{seq}} \mathbf{U}^* \times_2^{\text{seq}} \mathbf{U}^* + \mathcal{P}_m(\Delta_{\mathbf{T}}^{t,(m)}) \times_1^{\text{seq}} \mathbf{U}^* \times_2^{\text{seq}} \mathbf{U}^* \right), \quad (134)$$

$$\Delta_s^{t,(m)} := \mathbf{u}_s^{t,(m)} - \mathbf{u}_s^*, \quad (135)$$

for each $s \in [r]$.

Apply the triangle inequality to yield

$$\|(\mathbf{U}^{t+1,(m)} - \mathbf{U}^*)_{m,:}\|_2 \leq \underbrace{\|(\widehat{\mathbf{U}}^{t+1,(m)} - \mathbf{U}^*)_{m,:}\|_2}_{=:\alpha_1} + \underbrace{\|(\mathbf{U}^{t+1,(m)} - \widehat{\mathbf{U}}^{t+1,(m)})_{m,:}\|_2}_{=:\alpha_2},$$

leaving us with two terms to deal with. As it turns out, we will show that α_1 is the dominant term and α_2 is negligible. To simplify presentation, we shall assume that $\{E_{i,j,k}\}_{i,j,k \in [d]}$ (resp. $\{\chi_{i,j,k}\}_{i,j,k \in [d]}$) are independent random variables.

- Regarding α_1 , the definition of $\mathcal{P}_{\Omega_{-m}}$ allows us to derive

$$\begin{aligned}
(\widehat{\mathbf{U}}^{t+1,(m)} - \mathbf{U}^*)_{m,:} &= (\mathbf{U}^{t,(m)} - \mathbf{U}^* - \eta \Delta_{\mathbf{T}}^{t,(m)} \times_1^{\text{seq}} \mathbf{U}^* \times_2^{\text{seq}} \mathbf{U}^*)_{m,:} \\
&= (\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:} - \eta (\Delta_{\mathbf{T}}^{t,(m)} \times_1^{\text{seq}} \mathbf{U}^* \times_2^{\text{seq}} \mathbf{U}^*)_{m,:}.
\end{aligned}$$

We can express $\Delta_{\mathbf{T}}^{t,(m)} = \sum_{s \in [r]} (\Delta_s^{t,(m)} + \mathbf{u}_s^*)^{\otimes 3} - \mathbf{u}_s^{*\otimes 3}$ and compute that

$$\begin{aligned}
&(\Delta_{\mathbf{T}}^{t,(m)} \times_1 \mathbf{u}_s^* \times_2 \mathbf{u}_s^*)_m \\
&= \left(\|\mathbf{u}_s^*\|_2^2 + \langle \Delta_s^{t,(m)}, \mathbf{u}_s^* \rangle \right)^2 (\Delta_s^{t,(m)})_m + \left(2\|\mathbf{u}_s^*\|_2^2 \langle \Delta_s^{t,(m)}, \mathbf{u}_s^* \rangle + \langle \Delta_s^{t,(m)}, \mathbf{u}_s^* \rangle^2 \right) (\mathbf{u}_s^*)_m \\
&\quad + \sum_{i:i \neq s} \left(\langle \mathbf{u}_i^*, \mathbf{u}_s^* \rangle + \langle \Delta_i^{t,(m)}, \mathbf{u}_s^* \rangle \right)^2 (\Delta_i^{t,(m)})_m \\
&\quad + \sum_{i:i \neq s} \left(2\langle \mathbf{u}_i^*, \mathbf{u}_s^* \rangle \langle \Delta_i^{t,(m)}, \mathbf{u}_s^* \rangle + \langle \Delta_i^{t,(m)}, \mathbf{u}_s^* \rangle^2 \right) (\mathbf{u}_i^*)_m
\end{aligned} \quad (136)$$

for each $s \in [r]$. This further indicates that

$$\begin{aligned}
\left\| (\widehat{\mathbf{U}}^{t+1,(m)} - \mathbf{U}^*)_{m,:} \right\|_2 &\leq \underbrace{\left\| \sum_{s \in [r]} \left(1 - \eta \left(\|\mathbf{u}_s^*\|_2^2 + \langle \Delta_s^{t,(m)}, \mathbf{u}_s^* \rangle \right)^2 \right) (\Delta_s^{t,(m)})_m \mathbf{e}_s^\top \right\|_2}_{=:\beta_1} \\
&+ \eta \underbrace{\left\| \sum_{s \in [r]} \left(2 \|\mathbf{u}_s^*\|_2^2 \langle \Delta_s^{t,(m)}, \mathbf{u}_s^* \rangle + \langle \Delta_s^{t,(m)}, \mathbf{u}_s^* \rangle^2 \right) (\mathbf{u}_s^*)_m \mathbf{e}_s^\top \right\|_2}_{=:\beta_2} \\
&+ \eta \underbrace{\left\| \sum_{s \in [r]} \sum_{i:i \neq s} \left(\langle \mathbf{u}_i^*, \mathbf{u}_s^* \rangle + \langle \Delta_i^{t,(m)}, \mathbf{u}_s^* \rangle \right)^2 (\Delta_i^{t,(m)})_m \mathbf{e}_s^\top \right\|_2}_{=:\beta_3} \\
&+ \eta \underbrace{\left\| \sum_{s \in [r]} \sum_{i:i \neq s} \left(2 \langle \mathbf{u}_i^*, \mathbf{u}_s^* \rangle \langle \Delta_i^{t,(m)}, \mathbf{u}_s^* \rangle + \langle \Delta_i^{t,(m)}, \mathbf{u}_s^* \rangle^2 \right) (\mathbf{u}_i^*)_m \mathbf{e}_s^\top \right\|_2}_{=:\beta_4}.
\end{aligned}$$

In what follows, we will control the four terms separately.

– For β_1 , by (119), we use Cauchy-Schwarz to show that

$$\left(\|\mathbf{u}_s^*\|_2^2 + \langle \Delta_s^{t,(m)}, \mathbf{u}_s^* \rangle \right)^2 \geq \left(\|\mathbf{u}_s^*\|_2^2 - \|\Delta_s^{t,(m)}\|_2 \|\mathbf{u}_s^*\|_2 \right)^2 \geq \frac{2}{3} \|\mathbf{u}_s^*\|_2^4 \geq \frac{2}{3} \lambda_{\min}^{*4/3}$$

for each $s \in [r]$. It follows that

$$\beta_1 \leq \left(1 - \frac{2}{3} \lambda_{\min}^{*4/3} \eta \right) \left\| (\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:} \right\|_2. \quad (137)$$

– Regarding β_2 , by (119), we apply the Cauchy-Schwarz inequality again to get that: for each $s \in [r]$,

$$\left| \left(2 \|\mathbf{u}_s^*\|_2^2 \langle \Delta_s^{t,(m)}, \mathbf{u}_s^* \rangle + \langle \Delta_s^{t,(m)}, \mathbf{u}_s^* \rangle^2 \right) (\mathbf{u}_s^*)_m \right| \leq 3 \|\mathbf{u}_s^*\|_2^3 \|\mathbf{u}_s^*\|_\infty \|\Delta_s^{t,(m)}\|_2 \leq 3 \sqrt{\frac{\mu}{d}} \lambda_{\max}^{*4/3} \|\Delta_s^{t,(m)}\|_2,$$

Together with the assumption that $\kappa \asymp 1$, this implies that

$$\beta_2 \leq 3\eta \lambda_{\min}^{*4/3} \sqrt{\frac{\mu}{d}} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}}. \quad (138)$$

– With regards to β_3 , we show that for each $s \in [r]$:

$$\begin{aligned}
&\left| \sum_{i:i \neq s} \left(\langle \mathbf{u}_i^*, \mathbf{u}_s^* \rangle + \langle \Delta_i^{t,(m)}, \mathbf{u}_s^* \rangle \right)^2 (\Delta_i^{t,(m)})_m \right| \\
&\lesssim \max_{i:i \neq s} |\langle \mathbf{u}_i^*, \mathbf{u}_s^* \rangle|^2 \sum_{i:i \neq s} |(\Delta_i^{t,(m)})_m| + \|\mathbf{u}_s^*\|_2^2 \sum_{i:i \neq s} \|\Delta_i^{t,(m)}\|_2^2 |(\Delta_i^{t,(m)})_m| \\
&\leq \lambda_{\max}^{*4/3} \frac{\mu \sqrt{r}}{d} \left\| (\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:} \right\|_2 + \lambda_{\max}^{*2/3} \max_{i:i \neq s} \|\Delta_i^{t,(m)}\|_2 \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}} \left\| (\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:} \right\|_2 \\
&\leq \lambda_{\max}^{*4/3} \left(\frac{\mu \sqrt{r}}{d} + o(1/r) \right) \left\| (\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:} \right\|_2 \ll \frac{\lambda_{\max}^{*4/3}}{\sqrt{r}} \left\| (\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:} \right\|_2,
\end{aligned}$$

where the last line follows from (119) that $\max_{i:i \neq s} \|\Delta_i^{t,(m)}\|_2 \leq \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}} \ll \lambda_{\max}^{*1/3} / \sqrt{r}$ and the low rank condition $r \ll \sqrt{d/\mu}$. Summing over $s \in [r]$, we get

$$\beta_3 \ll \eta \lambda_{\min}^{*4/3} \left\| (\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:} \right\|_2 \quad (139)$$

under the condition $\kappa \asymp 1$.

– Turning attention to β_4 , we observe that for each $s \in [r]$,

$$\begin{aligned}
& \left| \sum_{i:i \neq s} \left(2 \langle \mathbf{u}_i^*, \mathbf{u}_s^* \rangle \langle \Delta_i^{t,(m)}, \mathbf{u}_s^* \rangle + \langle \Delta_i^{t,(m)}, \mathbf{u}_s^* \rangle^2 \right) (\mathbf{u}_i^*)_m \right| \\
& \lesssim \max_{i:i \neq s} |\langle \mathbf{u}_i^*, \mathbf{u}_s^* \rangle| \|\mathbf{u}_s^*\|_2 \sum_{i:i \neq s} \|\Delta_i^{t,(m)}\|_2 |(\mathbf{u}_i^*)_m| + \|\mathbf{u}_s^*\|_2^2 \sum_{i:i \neq s} \|\Delta_i^{t,(m)}\|_2^2 |(\mathbf{u}_i^*)_m| \\
& \leq \lambda_{\max}^* \sqrt{\frac{\mu}{d}} \|\mathbf{U}^*\|_{2,\infty} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\text{F}} + \lambda_{\max}^{*2/3} \max_{i:i \neq s} \|\Delta_i^{t,(m)}\|_2 \|\mathbf{U}^*\|_{2,\infty} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\text{F}} \\
& \lesssim \lambda_{\max}^* \sqrt{\frac{\mu r}{d}} \left(\sqrt{\frac{\mu}{d}} \lambda_{\max}^{*1/3} + o(1/\sqrt{r}) \right) \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\text{F}} \ll \lambda_{\max}^{*4/3} \sqrt{\frac{\mu}{d}} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\text{F}},
\end{aligned}$$

where the last line follows from (119) and the rank assumption $r \ll \sqrt{d/\mu}$. As a consequence,

$$\beta_4 \ll \eta \lambda_{\max}^{*4/3} \sqrt{\frac{\mu r}{d}} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\text{F}}. \quad (140)$$

– Therefore, we have

$$\left\| (\widehat{\mathbf{U}}^{t+1,(m)} - \mathbf{U}^*)_{m,:} \right\|_2 \leq \left(1 - \frac{1}{3} \lambda_{\min}^{*4/3} \eta \right) \left\| (\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:} \right\|_2 + 4\eta \lambda_{\max}^{*4/3} \sqrt{\frac{\mu r}{d}} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\text{F}}. \quad (141)$$

- With regards to α_2 , it follows from the definition (46) and (134) that

$$\begin{aligned}
& \mathbf{U}^{t+1,(m)} - \widehat{\mathbf{U}}^{t+1,(m)} \\
& = -\eta \left((p^{-1} \mathcal{P}_{\Omega_{-m}} + \mathcal{P}_m) (\Delta_{\mathbf{T}}^{t,(m)}) \times_1^{\text{seq}} \mathbf{U}^{t,(m)} \times_2^{\text{seq}} \mathbf{U}^{t,(m)} - (p^{-1} \mathcal{P}_{\Omega_{-m}} + \mathcal{P}_m) (\Delta_{\mathbf{T}}^{t,(m)}) \times_1^{\text{seq}} \mathbf{U}^* \times_2^{\text{seq}} \mathbf{U}^* \right) \\
& \quad + \frac{\eta}{p} \left(\mathcal{P}_{\Omega_{-m}}(\mathbf{E}) \times_1^{\text{seq}} \mathbf{U}^{t,(m)} \times_2^{\text{seq}} \mathbf{U}^{t,(m)} - \mathcal{P}_{\Omega_{-m}}(\mathbf{E}) \times_1^{\text{seq}} \mathbf{U}^* \times_2^{\text{seq}} \mathbf{U}^* \right).
\end{aligned} \quad (142)$$

Recall the definition of $\mathcal{P}_{\Omega_{-m}}$ and \mathcal{P}_m . For the m -th row, we have

$$(\mathbf{U}^{t+1,(m)} - \widehat{\mathbf{U}}^{t+1,(m)})_{m,:} = -\eta (\Delta_{\mathbf{T}}^{t,(m)} \times_1^{\text{seq}} \mathbf{U}^{t,(m)} \times_2^{\text{seq}} \mathbf{U}^{t,(m)} - \Delta_{\mathbf{T}}^{t,(m)} \times_1^{\text{seq}} \mathbf{U}^* \times_2^{\text{seq}} \mathbf{U}^*)_{m,:}$$

From the triangle inequality, we can further decompose

$$\begin{aligned}
& \left\| (\Delta_{\mathbf{T}}^{t,(m)} \times_1^{\text{seq}} \mathbf{U}^{t,(m)} \times_2^{\text{seq}} \mathbf{U}^{t,(m)} - \Delta_{\mathbf{T}}^{t,(m)} \times_1^{\text{seq}} \mathbf{U}^* \times_2^{\text{seq}} \mathbf{U}^*)_{m,:} \right\|_2 \\
& \leq \underbrace{\left\| (\Delta_{\mathbf{T}}^{t,(m)} \times_1^{\text{seq}} (\mathbf{U}^{t,(m)} - \mathbf{U}^*) \times_2^{\text{seq}} \mathbf{U}^*)_{m,:} \right\|_2}_{=:\gamma_1} + \underbrace{\left\| (\Delta_{\mathbf{T}}^{t,(m)} \times_1^{\text{seq}} \mathbf{U}^* \times_2^{\text{seq}} (\mathbf{U}^{t,(m)} - \mathbf{U}^*))_{m,:} \right\|_2}_{=:\gamma_2} \\
& \quad + \underbrace{\left\| (\Delta_{\mathbf{T}}^{t,(m)} \times_1^{\text{seq}} (\mathbf{U}^{t,(m)} - \mathbf{U}^*) \times_2^{\text{seq}} (\mathbf{U}^{t,(m)} - \mathbf{U}^*))_{m,:} \right\|_2}_{=:\gamma_3}.
\end{aligned}$$

Let us consider γ_1 first. It is straightforward to calculate that

$$\begin{aligned}
& (\Delta_{\mathbf{T}}^{t,(m)} \times_1 \Delta_s^{t,(m)} \times_2 \mathbf{u}_s^*)_m = \sum_{i \in [r]} \langle \mathbf{u}_i^{t,(m)}, \Delta_s^{t,(m)} \rangle \langle \mathbf{u}_i^{t,(m)}, \mathbf{u}_s^* \rangle (\Delta_i^{t,(m)})_m \\
& \quad + \sum_{i \in [r]} \left(\langle \Delta_i^{t,(m)}, \Delta_s^{t,(m)} \rangle \langle \mathbf{u}_i^{t,(m)}, \mathbf{u}_s^{t,(m)} \rangle + \langle \mathbf{u}_i^*, \Delta_s^{t,(m)} \rangle \langle \Delta_i^{t,(m)}, \mathbf{u}_s^* \rangle + \langle \Delta_i^{t,(m)}, \Delta_s^{t,(m)} \rangle \langle \Delta_i^{t,(m)}, \mathbf{u}_s^{t,(m)} \rangle \right) (\mathbf{u}_i^*)_m.
\end{aligned}$$

for each $s \in [r]$. From (119), we use the triangle inequality and the Cauchy-Schwarz inequality to upper bound

$$\left| (\Delta_{\mathbf{T}}^{t,(m)} \times_1 \Delta_s^{t,(m)} \times_2 \mathbf{u}_s^*)_m \right| \lesssim \|\Delta_s^{t,(m)}\|_2 \|\mathbf{u}_s^*\|_2 \sum_{i \in [r]} \|\mathbf{u}_i^{t,(m)}\|_2^2 |(\Delta_i^{t,(m)})_m|$$

$$\begin{aligned}
& + \|\Delta_s^{t,(m)}\|_2 (\|\mathbf{u}_s^*\|_2 + \|\mathbf{u}_s^{t,(m)}\|_2) \sum_{i \in [r]} \|\Delta_i^{t,(m)}\|_2 (\|\mathbf{u}_i^{t,(m)}\|_2 + \|\mathbf{u}_i^*\|_2 + \|\Delta_i^{t,(m)}\|_2) |(\mathbf{u}_i^*)_m| \\
& \lesssim \lambda_{\max}^{*2/3} \|\Delta_s^{t,(m)}\|_2 \sum_{i \in [r]} \|\mathbf{u}_i^{t,(m)}\|_2 \left| (\Delta_i^{t,(m)})_m \right| + \lambda_{\max}^{*2/3} \sqrt{\frac{\mu}{d}} \|\Delta_s^{t,(m)}\|_2 \sum_{i \in [r]} \|\Delta_i^{t,(m)}\|_2 \|\mathbf{u}_i^*\|_2 \\
& \lesssim \lambda_{\max}^{*2/3} \|\Delta_s^{t,(m)}\|_2 \|\mathbf{U}^{t,(m)}\|_{\mathbb{F}} \|(\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:}\|_2 + \lambda_{\max}^{*2/3} \sqrt{\frac{\mu}{d}} \|\Delta_s^{t,(m)}\|_2 \|\mathbf{U}^*\|_{\mathbb{F}} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}} \\
& \lesssim \lambda_{\max}^* \sqrt{r} \|\Delta_s^{t,(m)}\|_2 \|(\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:}\|_2 + \lambda_{\max}^* \sqrt{\frac{\mu r}{d}} \|\Delta_s^{t,(m)}\|_2 \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}}
\end{aligned}$$

We then sum over $s \in [r]$ to find that

$$\begin{aligned}
\left\| \left(\Delta_{\mathbf{T}}^{t,(m)} \times_1^{\text{seq}} (\mathbf{U}^{t,(m)} - \mathbf{U}^*) \times_2^{\text{seq}} \mathbf{U}^* \right)_{m,:} \right\|_2 & \lesssim \lambda_{\max}^* \sqrt{r} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}} \|(\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:}\|_2 \\
& + \lambda_{\max}^* \sqrt{\frac{\mu r}{d}} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}}^2, \tag{143}
\end{aligned}$$

Moreover, it is easy to see that the upper bound also holds for γ_2 . As for γ_3 , we can express

$$\begin{aligned}
(\Delta_{\mathbf{T}}^{t,(m)} \times_1 \Delta_s^{t,(m)} \times_2 \Delta_s^{t,(m)})_m & = \sum_{i \in [r]} \left(2 \langle \mathbf{u}_i^*, \Delta_s^{t,(m)} \rangle \langle \Delta_i^{t,(m)}, \Delta_s^{t,(m)} \rangle + \langle \Delta_i^{t,(m)}, \Delta_s^{t,(m)} \rangle^2 \right) (\mathbf{u}_i^*)_m \\
& + \sum_{i \in [r]} \langle \mathbf{u}_i^{t,(m)}, \Delta_s^{t,(m)} \rangle^2 (\Delta_i^{t,(m)})_m. \tag{144}
\end{aligned}$$

Similarly, we combine (119) with the triangle inequality and the Cauchy-Schwarz inequality to bound

$$\begin{aligned}
& \left| \left(\Delta_{\mathbf{T}}^{t,(m)} \times_1 \Delta_s^{t,(m)} \times_2 \Delta_s^{t,(m)} \right)_m \right| \\
& \lesssim \|\Delta_s^{t,(m)}\|_2^2 \sum_{i \in [r]} \|\Delta_i^{t,(m)}\|_2 \|\mathbf{u}_i^*\|_2 |(\mathbf{u}_i^*)_m| + \|\Delta_s^{t,(m)}\|_2^2 \sum_{i \in [r]} \|\mathbf{u}_i^{t,(m)}\|_2 \left| (\Delta_i^{t,(m)})_m \right| \\
& \lesssim \lambda_{\max}^{*1/3} \|\Delta_s^{t,(m)}\|_2^2 \sum_{i \in [r]} \|\Delta_i^{t,(m)}\|_2 |(\mathbf{u}_i^*)_m| + \lambda_{\max}^{*1/3} \|\Delta_s^{t,(m)}\|_2^2 \sum_{i \in [r]} \|\mathbf{u}_i^{t,(m)}\|_2 \left| (\Delta_i^{t,(m)})_m \right| \\
& \leq \lambda_{\max}^{*1/3} \|\Delta_s^{t,(m)}\|_2^2 \|\mathbf{U}^*\|_{2,\infty} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}} + \lambda_{\max}^{*1/3} \|\Delta_s^{t,(m)}\|_2^2 \|\mathbf{U}^{t,(m)}\|_{\mathbb{F}} \|(\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:}\|_2 \\
& \leq \lambda_{\max}^{*1/3} \|\mathbf{U}^*\|_{2,\infty} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}}^2 \|\Delta_s^{t,(m)}\|_2 \\
& + \lambda_{\max}^{*1/3} \|\mathbf{U}^*\|_{\mathbb{F}} \|(\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:}\|_2 \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}} \|\Delta_s^{t,(m)}\|_2.
\end{aligned}$$

Sum over $s \in [r]$ to obtain

$$\begin{aligned}
& \left\| \left(\Delta_{\mathbf{T}}^{t,(m)} \times_1^{\text{seq}} (\mathbf{U}^{t,(m)} - \mathbf{U}^*) \times_2^{\text{seq}} (\mathbf{U}^{t,(m)} - \mathbf{U}^*) \right)_{m,:} \right\|_2 \\
& \leq \lambda_{\max}^{*1/3} \sqrt{r} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}}^2 \|(\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:}\|_2 + \lambda_{\max}^{*1/3} \sqrt{\frac{\mu r}{d}} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}}^3.
\end{aligned}$$

Since $\|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}} \ll \lambda_{\max}^{*1/3} / \sqrt{r} \leq 1$ by (119), combined with (143), we find that

$$\begin{aligned}
& \left\| \left(\Delta_{\mathbf{T}}^{t,(m)} \times_1^{\text{seq}} \mathbf{U}^{t,(m)} \times_2^{\text{seq}} \mathbf{U}^{t,(m)} - \Delta_{\mathbf{T}}^{t,(m)} \times_1^{\text{seq}} \mathbf{U}^* \times_2^{\text{seq}} \mathbf{U}^* \right)_{m,:} \right\|_2 \\
& \lesssim \lambda_{\max}^{*1/3} \sqrt{r} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}} \|(\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:}\|_2 + \lambda_{\max}^{*1/3} \sqrt{\frac{\mu r}{d}} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}}^2 \\
& \ll \lambda_{\min}^{*4/3} \|(\mathbf{U}^{t,(m)} - \mathbf{U}^*)_{m,:}\|_2 + \lambda_{\min}^{*4/3} \sqrt{\frac{\mu r}{d}} \|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\mathbb{F}}. \tag{145}
\end{aligned}$$

where the last inequality follows from the assumption that $\kappa \asymp 1$.

- Putting (141) and (145) together, we reach the conclusion from (48) and the condition $r \ll \sqrt{d/\mu}$ that,

$$\begin{aligned}
\|(\mathbf{U}^{t+1,(m)} - \mathbf{U}^*)_{m,:}\|_2 &\leq \left(1 - \frac{1}{4}\lambda_{\min}^{*4/3}\eta\right) \|(\mathbf{U}^{t+1} - \mathbf{U}^*)_{m,:}\|_2 + 5\eta\lambda_{\min}^{*4/3}\sqrt{\frac{\mu r}{d}}\|\mathbf{U}^{t,(m)} - \mathbf{U}^*\|_{\text{F}} \\
&\leq \left(1 - \frac{1}{4}\lambda_{\min}^{*4/3}\eta\right) \left(C_7\rho^t\mathcal{E}_{\text{local}} + C_8\frac{\sigma}{\lambda_{\min}^*}\sqrt{\frac{d\log d}{p}}\right)\|\mathbf{U}^*\|_{2,\infty} \\
&\quad + 5\eta\lambda_{\min}^{*4/3}r\sqrt{\frac{\mu}{d}}\left(2C_1\rho^t\mathcal{E}_{\text{local}} + 2C_2\frac{\sigma}{\lambda_{\min}^*}\sqrt{\frac{d\log d}{p}}\right)\|\mathbf{U}^*\|_{2,\infty} \\
&\leq \left(\left(1 - \frac{1}{4}\lambda_{\min}^{*4/3}\eta\right)C_7 + o(1)C_1\lambda_{\min}^{*4/3}\eta\right)\rho^t\mathcal{E}_{\text{local}}\|\mathbf{U}^*\|_{2,\infty} \\
&\quad + \left(\left(1 - \lambda_{\min}^{*4/3}\eta\right)C_8 + o(1)C_2\lambda_{\min}^{*4/3}\eta\right)\frac{\sigma}{\lambda_{\min}^*}\sqrt{\frac{d\log d}{p}}\|\mathbf{U}^*\|_{2,\infty} \\
&\leq \left(C_7\rho^{t+1}\mathcal{E}_{\text{local}} + C_8\frac{\sigma}{\lambda_{\min}^*}\sqrt{\frac{d\log d}{p}}\right)\|\mathbf{U}^*\|_{2,\infty},
\end{aligned}$$

provided that $0 < \eta \leq \lambda_{\min}^{*4/3}/(32\lambda_{\max}^{*8/3})$, $1 - \lambda_{\min}^{*4/3}\eta/5 \leq \rho < 1$, C_7, C_8 are sufficiently large.

Recognizing that the above bound holds for any $1 \leq m \leq d$, we conclude the proof.

A.6 Proof of Lemma 5.6

It is easy to see that

$$\|(\mathbf{U}^{t+1} - \mathbf{U}^*)_{m,:}\|_2 \leq \|\mathbf{U}^{t+1} - \mathbf{U}^{t+1,(m)}\|_{\text{F}} + \|(\mathbf{U}^{t+1,(m)} - \mathbf{U}^*)_{m,:}\|_2. \quad (146)$$

Combining Lemma 5.4 and Lemma 5.5, we conclude that with probability at least $1 - O(d^{-10})$,

$$\begin{aligned}
\|(\mathbf{U}^{t+1} - \mathbf{U}^*)_{m,:}\|_2 &\leq \left(C_5\rho^t\mathcal{E}_{\text{local}} + C_6\frac{\sigma}{\lambda_{\min}^*}\sqrt{\frac{d\log d}{p}}\right)\|\mathbf{U}^*\|_{2,\infty} \\
&\quad + \left(C_7\rho^t\mathcal{E}_{\text{local}} + C_8\frac{\sigma}{\lambda_{\min}^*}\sqrt{\frac{d\log d}{p}}\right)\|\mathbf{U}^*\|_{2,\infty} \\
&\leq \left(C_3\rho^t\mathcal{E}_{\text{local}} + C_4\frac{\sigma}{\lambda_{\min}^*}\sqrt{\frac{d\log d}{p}}\right)\|\mathbf{U}^*\|_{2,\infty},
\end{aligned}$$

with the proviso that $C_3/(C_5 + C_7)$ and $C_4/(C_6 + C_8)$ are both sufficiently large.

B Proofs for retrieving tensor components

B.1 Proof of Lemma 5.12

We shall often operate upon the event where the claims in Lemma 5.7 hold, which happens with very high probability (i.e. at least $1 - O(d^{-10})$). Recall the definition of $\gamma^{*\tau}$ in (72). Since $\boldsymbol{\theta}^\tau = \mathbf{U}\mathbf{U}^\top\mathbf{g}^\tau$, this allows us to write that: for each $1 \leq i \leq r$,

$$\gamma_i^{*\tau} = \lambda_i^* \langle \mathbf{U}\mathbf{U}^\top \bar{\mathbf{u}}_i^*, \mathbf{g}^\tau \rangle,$$

where we recall that $\lambda_i^* = \|\mathbf{u}_i^*\|_2^3$ and $\bar{\mathbf{u}}_i^* = \mathbf{u}_i^*/\|\mathbf{u}_i^*\|_2$. Given that \mathbf{g}^τ is a Gaussian vector independent of \mathbf{U} , we observe that $\gamma^{*\tau}$ is zero-mean Gaussian conditional on Ω and \mathbf{E} . In order to understand the order statistics associated with this vector, we first look at its covariance matrix.

Denote by Σ^τ the covariance matrix of $\gamma^{*\tau}$ (conditional on \mathbf{U}). Then we have

$$\Sigma_{i,i}^\tau = \lambda_i^{*2} \|\mathbf{U}\mathbf{U}^\top \bar{\mathbf{u}}_i^*\|_2^2 = \lambda_i^{*2} \|\mathcal{P}_{\mathbf{U}}(\bar{\mathbf{u}}_i^*)\|_2^2, \quad (147a)$$

$$\begin{aligned} \Sigma_{i,j}^\tau &= \lambda_i^* \lambda_j^* \langle \mathbf{U}\mathbf{U}^\top \bar{\mathbf{u}}_i^*, \mathbf{U}\mathbf{U}^\top \bar{\mathbf{u}}_j^* \rangle = \lambda_i^* \lambda_j^* \langle \bar{\mathbf{u}}_i^*, \mathcal{P}_{\mathbf{U}}(\bar{\mathbf{u}}_j^*) \rangle \\ &= \lambda_i^* \lambda_j^* \left(\langle \bar{\mathbf{u}}_i^*, \bar{\mathbf{u}}_j^* \rangle - \langle \bar{\mathbf{u}}_i^*, \mathcal{P}_{\mathbf{U}^\perp}(\bar{\mathbf{u}}_j^*) \rangle \right), \end{aligned} \quad (147b)$$

where we denote by $\mathcal{P}_{\mathbf{U}}(\mathbf{z}) = \mathbf{U}\mathbf{U}^\top \mathbf{z}$ and $\mathcal{P}_{\mathbf{U}^\perp}(\mathbf{z}) = (\mathbf{I} - \mathbf{U}\mathbf{U}^\top) \mathbf{z}$. In addition, since the unit vector $\bar{\mathbf{u}}_i^*$ lies in the span of the columns of $\mathbf{U}_{\text{orth}}^*$ (cf. (54)), it follows from Lemma D.6 that

$$\|\mathcal{P}_{\mathbf{U}}(\bar{\mathbf{u}}_i^*)\|_2 = \|\mathbf{U}\mathbf{U}^\top \bar{\mathbf{u}}_i^*\|_2 = \|\mathbf{U}\mathbf{R}(\mathbf{U}\mathbf{R})^\top \bar{\mathbf{u}}_i^*\|_2 \geq \sqrt{1 - \|\mathbf{U}\mathbf{R} - \mathbf{U}_{\text{orth}}^*\|_2^2}, \quad (148a)$$

$$\|\mathcal{P}_{\mathbf{U}^\perp}(\bar{\mathbf{u}}_i^*)\|_2 = \|(\mathbf{I} - \mathbf{U}\mathbf{U}^\top) \bar{\mathbf{u}}_i^*\|_2 = \|(\mathbf{I} - \mathbf{U}\mathbf{R}(\mathbf{U}\mathbf{R})^\top) \bar{\mathbf{u}}_i^*\|_2 \leq \|\mathbf{U}\mathbf{R} - \mathbf{U}_{\text{orth}}^*\|_2, \quad (148b)$$

where \mathbf{R} is a rotation matrix defined in (55). This together with (147a) gives

$$\lambda_i^{*2} (1 - \|\mathbf{U}\mathbf{R} - \mathbf{U}_{\text{orth}}^*\|_2^2) \leq \Sigma_{i,i}^\tau \leq \lambda_i^{*2}, \quad (149)$$

where we have also used the fact that $\|\mathcal{P}_{\mathbf{U}}(\bar{\mathbf{u}}_i^*)\|_2 \leq \|\bar{\mathbf{u}}_i^*\|_2 = 1$. Moreover, taking together (147a), (148b) and the incoherence condition, we see that for any $1 \leq i \neq j \leq r$,

$$|\Sigma_{i,j}^\tau| \leq \lambda_i^* \lambda_j^* \left\{ |\langle \bar{\mathbf{u}}_i^*, \bar{\mathbf{u}}_j^* \rangle| + \|\bar{\mathbf{u}}_i^*\|_2 \|\mathcal{P}_{\mathbf{U}^\perp}(\bar{\mathbf{u}}_j^*)\|_2 \right\} \leq \lambda_i^* \lambda_j^* \underbrace{(\sqrt{\mu/d} + \|\mathbf{U}\mathbf{R} - \mathbf{U}_{\text{orth}}^*\|_2)}_{=: \delta_1}, \quad (150)$$

which is expected to be small if δ_1 is small.

From our assumptions on the sample size, the rank and the condition number, we can invoke Lemma 5.7 to see that $\kappa r \delta_1 \ll 1$, where δ_1 is defined in (150) and $\kappa = \lambda_{\max}^* / \lambda_{\min}^*$. Thus, we can decompose Σ^τ into two components as follows

$$\Sigma^\tau = \underbrace{(1 - \kappa r \delta_1) \mathbf{D}^{*2}}_{=: \hat{\Sigma}^\tau} + \underbrace{\Sigma^\tau - (1 - \kappa r \delta_1) \mathbf{D}^{*2}}_{=: \check{\Sigma}^\tau},$$

where

$$\mathbf{D}^* := \text{diag}(\lambda_1^*, \dots, \lambda_r^*) \in \mathbb{R}^{r \times r}.$$

As it turns out, both $\hat{\Sigma}^\tau$ and $\check{\Sigma}^\tau$ are positive definite. Indeed, we first learn from (149) and (150) that: the i -th diagonal entry of $\check{\Sigma}^\tau$ obeys

$$\begin{aligned} \check{\Sigma}_{i,i}^\tau &\geq \lambda_i^{*2} (1 - \|\mathbf{U}\mathbf{R} - \mathbf{U}_{\text{orth}}^*\|_2^2) - (1 - \kappa r \delta_1) \lambda_i^{*2} \\ &= \lambda_i^{*2} (\kappa r \delta_1 - \|\mathbf{U}\mathbf{R} - \mathbf{U}_{\text{orth}}^*\|_2^2) \geq \lambda_i^{*2} (\kappa r \delta_1 - \delta_1^2) \\ &\stackrel{(i)}{>} \kappa \lambda_i^{*2} (r-1) \delta_1 \stackrel{(ii)}{\geq} \lambda_i^* \lambda_{\max}^* (r-1) \delta_1 \\ &\geq \sum_{j:j \neq i} |\check{\Sigma}_{i,j}^\tau|, \end{aligned}$$

where (i) holds since $\delta_1 < 1$ under our assumptions, (ii) follows since $\kappa \lambda_i^* \geq \kappa \lambda_{\min}^* = \lambda_{\max}^*$, and the last line makes use of (150). This implies that $\check{\Sigma}^\tau$ is diagonally dominant, and hence $\check{\Sigma}^\tau \succeq \mathbf{0}$. In conclusion, both $\hat{\Sigma}^\tau$ and $\check{\Sigma}^\tau$ are positive definite.

Let $\hat{\gamma}^{*\tau}$ and $\check{\gamma}^{*\tau}$ be independent zero-mean Gaussian random vectors with covariance matrices $\hat{\Sigma}^\tau$ and $\check{\Sigma}^\tau$, respectively. Clearly, the distribution of $\gamma^{*\tau}$ is identical to that of $\hat{\gamma}^{*\tau} + \check{\gamma}^{*\tau}$. Consequently, it allows us to look at the distributions of these two random vectors separately.

- In view of (149) and the fact $\kappa r \delta_1 < 1$, one has

$$\check{\Sigma}_{i,i}^\tau \leq \lambda_i^{*2} - (1 - \kappa r \delta_1) \lambda_i^{*2} = \kappa r \delta_1 \lambda_i^{*2}. \quad (151)$$

Thus, with probability at least $1 - O(d^{-10})$, we have

$$\|\check{\gamma}^{*\tau}\|_\infty \lesssim \mathbb{E} \left[\max_{1 \leq i \leq r} |\check{\gamma}_i^{*\tau}| \right] + \max_{1 \leq i \leq r} \sqrt{\text{Var}(\check{\gamma}_i^{*\tau}) \log d} \lesssim \lambda_{\max}^* \sqrt{\kappa r \delta_1 \log d} \ll \lambda_{\min}^* / (r \sqrt{\log d}), \quad (152)$$

where the last step arises from $\kappa^3 r^3 \delta_1 \log^2 d \ll 1$ under our sample size, noise and rank conditions. By the condition that $L \asymp r^{2\kappa^2} \log^{3/2} r$, $r \ll d$ and $\kappa \asymp 1$, we take a union bound over $\tau \in [L]$ to find that each entry of $\check{\gamma}^{*\tau}$ is fairly small for all $\tau \in [L]$. Another immediate consequence of (152) is that: for any fixed vector $\mathbf{v} \in \mathbb{R}^r$, with probability at least $1 - O(d^{-10})$, for all $\tau \in [L]$,

$$\|\check{\gamma}^{*\tau}\|_2 \leq \sqrt{r} \|\check{\gamma}^{*\tau}\|_\infty \ll \lambda_{\min}^*, \quad (153a)$$

$$|\langle \mathbf{v}, \check{\gamma}^{*\tau} \rangle| \leq \|\mathbf{v}\|_2 \|\check{\gamma}^{*\tau}\|_2 \ll \|\mathbf{v}\|_2 \lambda_{\min}^*. \quad (153b)$$

- We then turn attention to $\hat{\gamma}^{*\tau}$, which is composed of independent Gaussian random variables. Let us define $\hat{\Delta}^\tau := \hat{\gamma}_1^{*\tau} - \max_{1 < i \leq r} |\hat{\gamma}_i^{*\tau}|$ for each $1 \leq \tau \leq L$ and let $\hat{\Delta}_1^{(1)} \geq \hat{\Delta}_1^{(2)} \geq \dots \geq \hat{\Delta}_1^{(L)}$ denote the order statistics of $\{\hat{\Delta}_1^\tau\}_{\tau=1}^L$ in descending order. Fix any small constant $\delta > 0$. Invoke Lemma D.5 to demonstrate that: with probability greater than $1 - \delta/r$,

$$\hat{\Delta}_1^{(1)} \gtrsim \lambda_{\min}^*, \quad (154a)$$

$$\hat{\Delta}_1^{(1)} - \hat{\Delta}_1^{(2)} \gtrsim \frac{\lambda_{\min}^*}{r \sqrt{\log d}}, \quad (154b)$$

where we use the conditions that $L \asymp r^{2\kappa^2} \log^{3/2} r$, $r \ll d$ and $\kappa \asymp 1$. In addition, let $\hat{\gamma}_{\setminus 1}^{*\tau} := [\hat{\gamma}_2^{*\tau}, \dots, \hat{\gamma}_r^{*\tau}]^\top \in \mathbb{R}^{r-1}$. We know from standard Gaussian concentration inequalities and union bounds that for any fixed vector $\mathbf{v} \in \mathbb{R}^r$, with probability $1 - O(d^{-20})$, for all $\tau \in [L]$,

$$\hat{\gamma}_1^{*\tau} \lesssim (\sqrt{\log L} + \sqrt{\log d}) \lambda_{\max}^* \asymp \sqrt{\log d} \lambda_{\max}^*, \quad (155a)$$

$$\|\hat{\gamma}^{*\tau}\|_2 \leq \hat{\gamma}_1^{*\tau} + \|\hat{\gamma}_{\setminus 1}^{*\tau}\|_2 \lesssim (\sqrt{\log d} + \sqrt{r \log d}) \lambda_{\max}^* \lesssim \sqrt{r \log d} \lambda_{\max}^*, \quad (155b)$$

$$|\langle \mathbf{v}, \hat{\gamma}^{*\tau} \rangle| \lesssim \hat{\gamma}_1^{*\tau} \|\mathbf{v}\|_2 + |\langle \mathbf{v}_{\setminus 1}, \hat{\gamma}_{\setminus 1}^{*\tau} \rangle| \lesssim \|\mathbf{v}\|_2 \sqrt{\log d} \lambda_{\max}^*, \quad (155c)$$

where $\mathbf{v}_{\setminus 1} := [v_2, \dots, v_r] \in \mathbb{R}^{r-1}$.

- Putting (152) and (154) together and invoking the triangle inequality immediately establish (74a) and (74b). On the other hand, combining (152) with (155a) proves (75a); (153a) and (155b) taken collectively establish (75b), whereas (153b) and (155c) prove (75c).

B.2 Proof of Lemma 5.13

Recall that the vector of interest $\bar{\mathbf{u}}^\tau$ is the leading singular vector of \mathbf{M}^τ (as constructed in (70c)), where \mathbf{M}^τ satisfies

$$\begin{aligned} \mathbf{M}^\tau &= p^{-1} \mathbf{T} \times_3 \boldsymbol{\theta}^\tau = \mathbf{T}^* \times_3 \boldsymbol{\theta}^\tau + (p^{-1} \mathbf{T} - \mathbf{T}^*) \times_3 \boldsymbol{\theta}^\tau \\ &= \gamma_1^{*\tau} \bar{\mathbf{u}}_1^* \bar{\mathbf{u}}_1^{*\top} + \sum_{s:s \neq 1} \gamma_s^{*\tau} \bar{\mathbf{u}}_s^* \bar{\mathbf{u}}_s^{*\top} + (p^{-1} \mathbf{T} - \mathbf{T}^*) \times_3 \boldsymbol{\theta}^\tau \\ &= \gamma_1^{*\tau} \bar{\mathbf{u}}_1^* \bar{\mathbf{u}}_1^{*\top} + \underbrace{\sum_{s:s \neq 1} \gamma_s^{*\tau} (\mathbf{I} - \bar{\mathbf{u}}_1^* \bar{\mathbf{u}}_1^{*\top}) \bar{\mathbf{u}}_s^* \bar{\mathbf{u}}_s^{*\top} (\mathbf{I} - \bar{\mathbf{u}}_1^* \bar{\mathbf{u}}_1^{*\top})}_{=: \mathbf{M}^{*\tau}} \\ &\quad + \underbrace{\sum_{s:s \neq 1} \gamma_s^{*\tau} \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle (\bar{\mathbf{u}}_1^* \bar{\mathbf{u}}_s^{*\top} + \bar{\mathbf{u}}_s^* \bar{\mathbf{u}}_1^{*\top}) - \sum_{s:s \neq 1} \gamma_s^{*\tau} \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle^2 \bar{\mathbf{u}}_s^* \bar{\mathbf{u}}_s^{*\top}}_{=: \mathbf{C}^\tau} + \underbrace{(p^{-1} \mathbf{T} - \mathbf{T}^*) \times_3 \boldsymbol{\theta}^\tau}_{=: \mathbf{F}^\tau}, \quad (156) \end{aligned}$$

and $\gamma_i^{*\tau}$ ($1 \leq i \leq r$) is defined in (72).

In what follows, we shall view \mathbf{C}^τ and \mathbf{F}^τ as perturbation terms superimposed on $\mathbf{M}^{*\tau}$. Lemma B.1 below proves that their operator norms are all small under our sample size, noise and rank conditions, which enables to apply Wedin's theorem to justify the ℓ_2 proximity between $\bar{\mathbf{u}}^\tau$ and $\bar{\mathbf{u}}_1^*$.

Lemma B.1. *Instate the assumptions of Lemma 5.13. With probability at least $1 - O(d^{-10})$, one has*

$$\|\mathbf{F}^\tau\| \lesssim \frac{\sqrt{\mu r} \lambda_{\max}^* \log^3 d}{d^{3/2} p} + \frac{\mu r \lambda_{\max}^* \log^{5/2} d}{d \sqrt{p}} + \frac{\sigma \log^{7/2} d}{p} + \sigma \sqrt{\frac{rd \log^5 d}{p}}, \quad (157)$$

$$\|\mathbf{C}^\tau\| \lesssim \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^*, \quad (158)$$

$$\|\mathbf{F}^\tau \bar{\mathbf{u}}_1^*\|_2 \lesssim \frac{\mu r \lambda_{\max}^* \log d}{d \sqrt{p}} + \sigma \sqrt{\frac{rd \log^2 d}{p}}, \quad (159)$$

$$\|\mathbf{C}^\tau \bar{\mathbf{u}}_1^*\|_2 \lesssim \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^*. \quad (160)$$

Proof. See Appendix B.3. ■

As a consequence, recalling the definition of \mathcal{E}_{op} and $\mathcal{E}_{\text{proj}}$ in (81) and (79) respectively, one has

$$\begin{aligned} \|\mathbf{M}^\tau - \mathbf{M}^{*\tau}\| &\leq \|\mathbf{F}^\tau\| + \|\mathbf{C}^\tau\| \\ &\lesssim \underbrace{\frac{\sqrt{\mu r} \lambda_{\max}^* \log^3 d}{d^{3/2} p} + \frac{\mu r \lambda_{\max}^* \log^{5/2} d}{d \sqrt{p}} + \frac{\sigma \log^{7/2} d}{p} + \sigma \sqrt{\frac{rd \log^5 d}{p}} + \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^*}_{=\mathcal{E}_{\text{op}} \cdot \lambda_{\min}^*}, \end{aligned} \quad (161)$$

and

$$\|(\mathbf{M}^\tau - \mathbf{M}^{*\tau}) \bar{\mathbf{u}}_1^*\|_2 \leq \|\mathbf{F}^\tau \bar{\mathbf{u}}_1^*\|_2 + \|\mathbf{C}^\tau \bar{\mathbf{u}}_1^*\|_2 \lesssim \underbrace{\frac{\mu r \lambda_{\max}^* \log d}{d \sqrt{p}} + \sigma \sqrt{\frac{rd \log^2 d}{p}} + \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^*}_{=\mathcal{E}_{\text{proj}} \cdot \lambda_{\min}^*}. \quad (162)$$

It then follows from Weyl's inequality that

$$\max_{i \in [d]} |\sigma_i(\mathbf{M}^\tau) - \sigma_i(\mathbf{M}^{*\tau})| \leq \|\mathbf{M}^\tau - \mathbf{M}^{*\tau}\| \lesssim \mathcal{E}_{\text{op}} \cdot \lambda_{\min}^* \ll \lambda_{\min}^* \quad (163)$$

where $\sigma_i(\mathbf{Z})$ denotes the i -th largest singular value of a matrix \mathbf{Z} and we use the condition that $\mathcal{E}_{\text{op}} \ll 1$. All in all, these arguments justify that the spectrum of \mathbf{M}^τ is fairly close to that of $\mathbf{M}^{*\tau}$.

Next, we look at the gap between the two leading singular values of $\mathbf{M}^{*\tau}$. To begin with, it is self-evident from the definition of $\mathbf{M}^{*\tau}$ that: $\bar{\mathbf{u}}_1^*$ is the singular vector of $\mathbf{M}^{*\tau}$. In fact, we claim one further result, that is, $\bar{\mathbf{u}}_1^*$ is indeed the leading singular vector of $\mathbf{M}^{*\tau}$ whose singular value is given by $\sigma_1(\mathbf{M}^{*\tau}) = \gamma_1^{*\tau}$. Towards this end, let us define

$$\mathbf{U}_{\prec 1}^\tau := (\mathbf{I} - \bar{\mathbf{u}}_1^* \bar{\mathbf{u}}_1^{*\top}) \bar{\mathbf{U}}^* \in \mathbb{R}^{d \times (r-1)} \quad \text{and} \quad \boldsymbol{\gamma}_{\prec 1}^{*\tau} := [\gamma_2^{*\tau}, \dots, \gamma_r^{*\tau}]^\top \in \mathbb{R}^{r-1},$$

allowing us to write

$$\sum_{s: s \neq 1} \gamma_s^{*\tau} (\mathbf{I} - \bar{\mathbf{u}}_1^* \bar{\mathbf{u}}_1^{*\top}) \bar{\mathbf{u}}_s^* \bar{\mathbf{u}}_s^{*\top} (\mathbf{I} - \bar{\mathbf{u}}_1^* \bar{\mathbf{u}}_1^{*\top}) = \mathbf{U}_{\prec 1}^\tau \text{diag}(\boldsymbol{\gamma}_{\prec 1}^{*\tau}) \mathbf{U}_{\prec 1}^{\tau\top} =: \mathbf{M}_{\prec 1}^{*\tau}.$$

We note that from Lemma D.1, one has

$$\|\mathbf{U}_{\prec 1}^\tau\| = \|(\mathbf{I} - \bar{\mathbf{u}}_1^* \bar{\mathbf{u}}_1^{*\top}) \bar{\mathbf{U}}^*\| \leq \|\bar{\mathbf{U}}^*\| \leq \sqrt{1 + r \sqrt{\mu/d}}.$$

Let $|\gamma^{*\tau}|_{(1)} \geq \dots \geq |\gamma^{*\tau}|_{(r)}$ denote the absolute values of $\{\gamma_i^{*\tau}\}_{i=1}^r$ in descending order. This together with Lemma 5.12 implies that

$$\|\mathbf{M}_{\prec 1}^{*\tau}\| \leq |\gamma^{*\tau}|_{(2)} \|\mathbf{U}_{\prec 1}^\tau\|^2 \leq |\gamma^{*\tau}|_{(2)} (1 + r \sqrt{\mu/d})$$

$$\leq \gamma_1^{*\tau} - (\gamma_1^{*\tau} - |\gamma^{*\tau}|_{(2)}) + r\sqrt{\mu/d}\gamma_1^{*\tau} < \gamma_1^{*\tau},$$

as long as $\kappa r\sqrt{(\mu \log d)/d} \ll 1$. Given that $\bar{\mathbf{u}}_1^*$ is the singular vector of $\mathbf{M}^{*\tau}$ with singular value $\gamma_1^{*\tau}$, we can conclude that $\sigma_1(\mathbf{M}^{*\tau}) = \gamma_1^{*\tau}$. This also allows us to lower bound the gap between the two largest singular values $\mathbf{M}^{*\tau}$ as follows

$$\begin{aligned} \sigma_1(\mathbf{M}^{*\tau}) - \sigma_2(\mathbf{M}^{*\tau}) &\geq \gamma_1^{*\tau} - \|\mathbf{M}_{\leq 1}^{*\tau}\| \geq \gamma_1^{*\tau} - |\gamma^{*\tau}|_{(2)}(1 + r\sqrt{\mu/d}) \\ &\gtrsim \gamma_1^{*\tau} - |\gamma^{*\tau}|_{(2)} \gtrsim \lambda_{\min}^*, \end{aligned} \quad (164)$$

provided that $\kappa r\sqrt{(\mu \log d)/d} \ll 1$. We also know from (163) and (164) that

$$\sigma_1(\mathbf{M}^\tau) - \sigma_2(\mathbf{M}^\tau) \geq \sigma_1(\mathbf{M}^{*\tau}) - \sigma_2(\mathbf{M}^{*\tau}) - 2\|\mathbf{M}^\tau - \mathbf{M}^{*\tau}\| \gtrsim \lambda_{\min}^*.$$

Combined with (162) and Wedin's theorem, we conclude that

$$\|\bar{\mathbf{u}}^\tau - \bar{\mathbf{u}}_1^*\|_2 \leq \frac{\|(\mathbf{M}^\tau - \mathbf{M}^{*\tau})\bar{\mathbf{u}}_1^*\|_2}{\sigma_1(\mathbf{M}^\tau) - \sigma_2(\mathbf{M}^\tau) - \|\mathbf{M}^\tau - \mathbf{M}^{*\tau}\|} \lesssim \frac{\mu r \log d}{d\sqrt{p}} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{rd \log^2 d}{p}} + \sqrt{\frac{\mu r \log d}{d}}. \quad (165)$$

Here, we have made use of the fact that $\bar{\mathbf{u}}^\tau$ is the leading singular vector of \mathbf{M}^τ obeying $\langle \bar{\mathbf{u}}^\tau, \bar{\mathbf{u}}_1^* \rangle \geq 0$.

B.3 Proof of Lemma B.1

B.3.1 Controlling \mathbf{F}^τ

- We first consider the spectral norm of \mathbf{F}^τ . Recall the definition that $\boldsymbol{\theta}^\tau = \mathbf{U}\mathbf{U}^\top \mathbf{g}^\tau$. Let us define

$$\boldsymbol{\theta}^{*\tau} := \mathbf{U}_{\text{orth}}^* \mathbf{U}_{\text{orth}}^{*\top} \mathbf{g}^\tau$$

and decompose

$$\mathbf{F}^\tau = (p^{-1}\mathbf{T} - \mathbf{T}^*) \times_3 \boldsymbol{\theta}^\tau = \underbrace{(p^{-1}\mathbf{T} - \mathbf{T}^*) \times_3 \boldsymbol{\theta}^{*\tau}}_{=: \mathbf{X}} + \underbrace{(p^{-1}\mathbf{T} - \mathbf{T}^*) \times_3 (\boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{*\tau})}_{=: \mathbf{Y}}.$$

In the sequel, we shall control these two terms separately.

- To bound $\|\mathbf{X}\|$, observe that $\boldsymbol{\theta}^{*\tau}$ is independent of $p^{-1}\mathbf{T} - \mathbf{T}^*$. By Lemma D.4, one has

$$\|(p^{-1}\mathbf{T} - \mathbf{T}^*) \times_3 \boldsymbol{\theta}^{*\tau}\| \lesssim \|\boldsymbol{\theta}^{*\tau}\|_\infty \sqrt{\frac{\mu r \log d}{dp}} \lambda_{\max}^* + \|\boldsymbol{\theta}^{*\tau}\|_\infty \frac{\sigma \log^{5/2} d}{p} + \|\boldsymbol{\theta}^{*\tau}\|_2 \sigma \sqrt{\frac{d \log d}{p}}. \quad (166)$$

This suggests that we need to control the ℓ_∞ and ℓ_2 norms of $\boldsymbol{\theta}^{*\tau}$. Using standard results on Gaussian random vectors and Lemma D.1, we know that with probability at least $1 - O(d^{-20})$,

$$\|\boldsymbol{\theta}^{*\tau}\|_\infty = \|\mathbf{U}_{\text{orth}}^* \mathbf{U}_{\text{orth}}^{*\top} \mathbf{g}^\tau\|_\infty \lesssim \|\mathbf{U}_{\text{orth}}^*\|_{2,\infty} \sqrt{\log d} \leq \sqrt{\frac{\mu r \log d}{d}}, \quad (167)$$

$$\|\boldsymbol{\theta}^{*\tau}\|_2 = \|\mathbf{U}_{\text{orth}}^* \mathbf{U}_{\text{orth}}^{*\top} \mathbf{g}^\tau\|_2 \lesssim \|\mathbf{U}_{\text{orth}}^*\|_{\text{F}} \sqrt{\log d} = \sqrt{r \log d}. \quad (168)$$

Combining (166) with (167) and (168) reveals that with probability exceeding $1 - O(d^{-20})$,

$$\begin{aligned} \|(p^{-1}\mathbf{T} - \mathbf{T}^*) \times_3 \boldsymbol{\theta}^{*\tau}\| &\lesssim \frac{\mu r \lambda_{\max}^* \log d}{d\sqrt{p}} + \frac{\sigma}{p} \sqrt{\frac{\mu r \log^6 d}{d}} + \sigma \sqrt{\frac{rd \log^2 d}{p}} \\ &\asymp \frac{\mu r \lambda_{\max}^* \log d}{d\sqrt{p}} + \sigma \sqrt{\frac{rd \log^2 d}{p}}, \end{aligned} \quad (169)$$

where the last inequality holds as long as $p \gtrsim \mu d^{-2} \log^4 d$.

– Turning to \mathbf{Y} , we can simply upper bound

$$\|(p^{-1}\mathbf{T} - \mathbf{T}^*) \times_3 (\boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{*\tau})\| \leq \|p^{-1}\mathbf{T} - \mathbf{T}^*\| \|\boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{*\tau}\|_2.$$

Since $\text{rank}(\mathbf{U}\mathbf{U}^\top - \mathbf{U}_{\text{orth}}^* \mathbf{U}_{\text{orth}}^{*\top}) \leq 2r$, Lemma 5.7 and the standard result of Gaussian random vectors yields that: with probability at least $1 - O(d^{-12})$,

$$\begin{aligned} \|\boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{*\tau}\|_2 &= \|(\mathbf{U}\mathbf{U}^\top - \mathbf{U}_{\text{orth}}^* \mathbf{U}_{\text{orth}}^{*\top}) \mathbf{g}^\tau\|_2 \lesssim \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}_{\text{orth}}^* \mathbf{U}_{\text{orth}}^{*\top}\|_{\text{F}} \sqrt{\log d} \\ &\leq \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}_{\text{orth}}^* \mathbf{U}_{\text{orth}}^{*\top}\| \sqrt{2r \log d} \lesssim \mathcal{E}_{\text{se}} \sqrt{r \log d} \ll 1, \end{aligned} \quad (170)$$

where we recall the definition of \mathcal{E}_{se} in (57) and that $\mathcal{E}_{\text{se}} \ll 1/\sqrt{r \log d}$ by our conditions. Moreover, by Lemma D.2, we know that with probability exceeding $1 - O(d^{-10})$,

$$\begin{aligned} \|p^{-1}\mathbf{T} - \mathbf{T}^*\| &\leq \|p^{-1}\mathcal{P}_\Omega(\mathbf{T}^*) - \mathbf{T}^*\| + \|p^{-1}\mathcal{P}_\Omega(\mathbf{E})\| \\ &\lesssim \frac{\sqrt{\mu r} \lambda_{\max}^* \log^3 d}{d^{3/2} p} + \frac{\mu \sqrt{r} \lambda_{\max}^* \log^{5/2} d}{d \sqrt{p}} + \frac{\sigma \log^{7/2} d}{p} + \sigma \sqrt{\frac{d \log^5 d}{p}}. \end{aligned} \quad (171)$$

Combining (170) and (171), we find that

$$\|(p^{-1}\mathbf{T} - \mathbf{T}^*) \times_3 (\boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{*\tau})\| \lesssim \frac{\sqrt{\mu r} \lambda_{\max}^* \log^3 d}{d^{3/2} p} + \frac{\mu \sqrt{r} \lambda_{\max}^* \log^{5/2} d}{d \sqrt{p}} + \frac{\sigma \log^{7/2} d}{p} + \sigma \sqrt{\frac{d \log^5 d}{p}}. \quad (172)$$

Putting (169) and (172) together shows that

$$\|\mathbf{F}^\tau\| \lesssim \frac{\sqrt{\mu r} \lambda_{\max}^* \log^3 d}{d^{3/2} p} + \frac{\mu r \lambda_{\max}^* \log^{5/2} d}{d \sqrt{p}} + \frac{\sigma \log^{7/2} d}{p} + \sigma \sqrt{\frac{rd \log^5 d}{p}}.$$

- Next, we turn to $\|\mathbf{F}^\tau \bar{\mathbf{u}}_1^*\|_2$. By the definition of the operator norm, we know that

$$\|\mathbf{F}^\tau \bar{\mathbf{u}}_1^*\|_2 = \|(p^{-1}\mathbf{T} - \mathbf{T}^*) \times_2 \bar{\mathbf{u}}_1^* \times_3 \boldsymbol{\theta}^\tau\|_2 \leq \|(p^{-1}\mathbf{T} - \mathbf{T}^*) \times_2 \bar{\mathbf{u}}_1^*\| \|\boldsymbol{\theta}^\tau\|_2.$$

Applying Lemma D.4 again reveals that with probability at least $1 - O(d^{-12})$,

$$\begin{aligned} \|(p^{-1}\mathbf{T} - \mathbf{T}^*) \times_2 \bar{\mathbf{u}}_1^*\| &\lesssim \|\bar{\mathbf{u}}_1^*\|_\infty \sqrt{\frac{\mu r \log d}{dp}} \lambda_{\max}^* + \|\bar{\mathbf{u}}_1^*\|_\infty \frac{\sigma \log^{5/2} d}{p} + \|\bar{\mathbf{u}}_1^*\|_2 \sigma \sqrt{\frac{d \log d}{p}} \\ &\lesssim \frac{\mu \sqrt{r} \lambda_{\max}^* \sqrt{\log d}}{d \sqrt{p}} + \frac{\sigma \sqrt{\mu \log^5 d}}{\sqrt{d} p} + \sigma \sqrt{\frac{d \log d}{p}} \\ &\asymp \frac{\mu \sqrt{r} \lambda_{\max}^* \sqrt{\log d}}{d \sqrt{p}} + \sigma \sqrt{\frac{d \log d}{p}}, \end{aligned} \quad (173)$$

where the last step arises from the condition that $p \gtrsim \mu d^{-2} \log^4 d$. In addition, realizing that \mathbf{U} consists of eigenvectors, standard Gaussian random vectors results give that with probability at least $1 - O(d^{-12})$,

$$\|\boldsymbol{\theta}^\tau\|_2 = \|\mathbf{U}\mathbf{U}^\top \mathbf{g}^\tau\|_2 \lesssim \|\mathbf{U}\|_{\text{F}} \sqrt{\log d} = \sqrt{r \log d}. \quad (174)$$

Combining (173) and (174) shows that with probability exceeding $1 - O(d^{-10})$,

$$\|\mathbf{F}^\tau \bar{\mathbf{u}}_1^*\|_2 \lesssim \frac{\mu r \lambda_{\max}^* \log d}{d \sqrt{p}} + \sigma \sqrt{\frac{rd \log^2 d}{p}}.$$

B.3.2 Controlling C^τ

Recall the definition of C^τ in (156). We first consider the spectral norm of C^τ . It is straightforward to compute that

$$\begin{aligned} \|C^\tau\| &\leq 2 \left\| \sum_{s:s \neq 1} \gamma_s^{*\tau} \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle \bar{\mathbf{u}}_s^* \right\|_2 + \left| \sum_{s:s \neq 1} \gamma_s^{*\tau} \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle^2 \right| \\ &\lesssim \max_{s:s \neq 1} |\langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle| \|\bar{\mathbf{U}}^*\| \|\gamma^{*\tau}\|_2 + \max_{s:s \neq 1} \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle^2 \sqrt{r} \|\gamma^{*\tau}\|_2 \\ &\lesssim \left(\sqrt{\frac{\mu}{d}} + \frac{\mu\sqrt{r}}{d} \right) \sqrt{r \log d} \lambda_{\max}^* \asymp \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^* \end{aligned} \quad (175)$$

if $\mu r/d \lesssim 1$, where we recall that $\bar{\mathbf{U}}^* = [\bar{\mathbf{u}}_1^*, \dots, \bar{\mathbf{u}}_r^*]$. Here, the last line holds owing to (8c), Lemma D.1 (which justifies that $\|\bar{\mathbf{U}}^*\| \lesssim 1$ if $r\sqrt{\mu/d} \leq 1$) and Lemma 5.12.

The claim (160) arises from the definition of the spectral norm that $\|C^\tau \bar{\mathbf{u}}_1^*\|_2 \leq \|C^\tau\|$.

B.4 Proof of Lemma 5.14

Let us fix an arbitrary $m \in [d]$. We remind the readers of several definitions: (1) $\gamma_1^{*\tau}$: see (72); (2) $M^{*\tau}$: see (156); and (3) $M^{\tau,(m)}$: see (70d).

Before continuing, we state two immediate facts. First, it has already been observed in Appendix B.2 that $\bar{\mathbf{u}}_1^*$ is a singular vector of $M^{*\tau}$ with singular value $\gamma_1^{*\tau}$, and hence

$$(\bar{\mathbf{u}}_1^*)_m = (\gamma_1^{*\tau})^{-1} M_{m,:}^{*\tau} \bar{\mathbf{u}}_1^*. \quad (176)$$

Here and throughout, $\mathbf{A}_{m,:}$ denotes the m -th row of a matrix \mathbf{A} . Second, $\bar{\mathbf{u}}^{\tau,(m)}$ is the top singular vector of $M^{\tau,(m)}$ such that $\langle \bar{\mathbf{u}}^{\tau,(m)}, \bar{\mathbf{u}}_1^* \rangle \geq 0$, and we denote by $\gamma_\tau^{(m)}$ the associated singular value. Recall our definition of $\nu^{\tau,(m)}$ in Algorithm 6. Similar to the case of ν^τ , we will show shortly in Lemma 5.16 that the global signs of $\nu^{\tau,(m)}$ and $\bar{\mathbf{u}}^{\tau,(m)}$ coincide, and hence

$$\nu^{\tau,(m)} = \bar{\mathbf{u}}^{\tau,(m)}. \quad (177)$$

As a result, the proof of this lemma boils down to showing that $\bar{\mathbf{u}}^{\tau,(m)}$ (and hence $\nu^{\tau,(m)}$) is sufficiently close to $\bar{\mathbf{u}}_1^*$ in the m -th entry. Towards this end, observe that

$$(\bar{\mathbf{u}}^{\tau,(m)})_m = (\gamma_\tau^{(m)})^{-1} M_{m,:}^{\tau,(m)} \bar{\mathbf{u}}^{\tau,(m)}. \quad (178)$$

The above two facts (176) and (178) together with the triangle inequality lead to

$$\begin{aligned} \left| (\bar{\mathbf{u}}^{\tau,(m)} - \bar{\mathbf{u}}_1^*)_m \right| &\leq \left| \left\{ (\gamma_\tau^{(m)})^{-1} - (\gamma_1^{*\tau})^{-1} \right\} M_{m,:}^{\tau,(m)} \bar{\mathbf{u}}^{\tau,(m)} \right| \\ &\quad + (\gamma_1^{*\tau})^{-1} \left| (M^{\tau,(m)} - M^{*\tau})_{m,:} \bar{\mathbf{u}}^{\tau,(m)} \right| \\ &\quad + (\gamma_1^{*\tau})^{-1} \left| M_{m,:}^{*\tau} (\bar{\mathbf{u}}^{\tau,(m)} - \bar{\mathbf{u}}_1^*) \right| \\ &\leq \underbrace{\left| (\gamma_\tau^{(m)})^{-1} - (\gamma_1^{*\tau})^{-1} \right| \|M_{m,:}^{\tau,(m)}\|_2 \|\bar{\mathbf{u}}^{\tau,(m)}\|_2}_{=:\alpha_1} \\ &\quad + \underbrace{(\gamma_1^{*\tau})^{-1} \|(M^{\tau,(m)} - M^{*\tau})_{m,:}\|_2 \|\bar{\mathbf{u}}^{\tau,(m)}\|_2}_{=:\alpha_2} \\ &\quad + \underbrace{(\gamma_1^{*\tau})^{-1} \|M_{m,:}^{*\tau}\|_2 \|\bar{\mathbf{u}}^{\tau,(m)} - \bar{\mathbf{u}}_1^*\|_2}_{=:\alpha_3}. \end{aligned} \quad (179)$$

Therefore, it suffices to upper bound the above three quantities separately.

B.4.1 Controlling α_3

The first step to bound α_3 (cf. (179)) is to control $\|\mathbf{M}_{m,:}^{*\tau}\|_2$. Towards this end, we first observe from the incoherence conditions that

$$\begin{aligned} \max_{s:s \neq 1} |(\bar{\mathbf{u}}_s^* - \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle \bar{\mathbf{u}}_1^*)_m| &\leq \max_{s:s \neq 1} \|\bar{\mathbf{u}}_s^*\|_\infty + \max_{s:s \neq 1} |\langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle| \|\bar{\mathbf{u}}_1^*\|_\infty \\ &\leq \max_{s:s \neq 1} \|\bar{\mathbf{u}}_s^*\|_\infty + \max_{s:s \neq 1} \|\bar{\mathbf{u}}_s^*\|_2 \|\bar{\mathbf{u}}_1^*\|_2 \|\bar{\mathbf{u}}_1^*\|_\infty \lesssim \sqrt{\frac{\mu}{d}}. \end{aligned} \quad (180)$$

When combined with the definition (156), this gives

$$\begin{aligned} \|\mathbf{M}_{m,:}^{*\tau}\|_2 &= \left\| \gamma_1^{*\tau} (\bar{\mathbf{u}}_1^*)_m \bar{\mathbf{u}}_1^{*\top} + \sum_{s:s \neq 1} \gamma_s^{*\tau} (\bar{\mathbf{u}}_s^* - \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle \bar{\mathbf{u}}_1^*)_m \bar{\mathbf{u}}_s^{*\top} (\mathbf{I} - \bar{\mathbf{u}}_1^* \bar{\mathbf{u}}_1^{*\top}) \right\|_2 \\ &\leq \gamma_1^{*\tau} \|\bar{\mathbf{u}}_1^*\|_\infty \|\bar{\mathbf{u}}_1^*\|_2 + \left\| \sum_{s:s \neq 1} \gamma_s^{*\tau} (\bar{\mathbf{u}}_s^* - \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle \bar{\mathbf{u}}_1^*)_m \bar{\mathbf{u}}_s^{*\top} \right\|_2 \\ &\leq \gamma_1^{*\tau} \|\bar{\mathbf{u}}_1^*\|_\infty + \max_{s:s \neq 1} |(\bar{\mathbf{u}}_s^* - \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle \bar{\mathbf{u}}_1^*)_m| \|\gamma^{*\tau}\|_2 \|\bar{\mathbf{U}}^*\| \\ &\stackrel{(i)}{\lesssim} \gamma_1^{*\tau} \sqrt{\frac{\mu}{d}} + \|\gamma^{*\tau}\|_2 \sqrt{\frac{\mu}{d}} \asymp \|\gamma^{*\tau}\|_2 \sqrt{\frac{\mu}{d}} \lesssim \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^*, \end{aligned} \quad (181)$$

where (i) arises from (180) and $\|\bar{\mathbf{U}}^*\| \lesssim 1$ if $r\sqrt{\mu/d} \ll 1$, and the last step comes from Lemma 5.12.

The second step is to upper bound $\|\bar{\mathbf{u}}^{\tau,(m)} - \bar{\mathbf{u}}_1^*\|_2$. Towards this, we resort to Wedin's theorem as follows

$$\|\bar{\mathbf{u}}^{\tau,(m)} - \bar{\mathbf{u}}_1^*\|_2 \leq \frac{\|(\mathbf{M}^{\tau,(m)} - \mathbf{M}^{*\tau}) \bar{\mathbf{u}}_1^*\|_2}{\sigma_1(\mathbf{M}^{*\tau}) - \sigma_2(\mathbf{M}^{*\tau}) - \|\mathbf{M}^{\tau,(m)} - \mathbf{M}^{*\tau}\|}, \quad (182)$$

where we rely on the fact that $\langle \bar{\mathbf{u}}^{\tau,(m)}, \bar{\mathbf{u}}_1^* \rangle \geq 0$. To complete this bound, we need to control $\mathbf{M}^{\tau,(m)} - \mathbf{M}^{*\tau}$. Before we move on, we find it helpful to introduce

$$\widehat{\mathbf{M}}^{\tau,(m)} := p^{-1} \mathbf{T} \times_3 \boldsymbol{\theta}^{\tau,(m)}. \quad (183)$$

Let $\widehat{\mathbf{u}}^{\tau,(m)}$ denote the top left singular vector of $\widehat{\mathbf{M}}^{\tau,(m)}$ such that

$$\langle \widehat{\mathbf{u}}^{\tau,(m)}, \bar{\mathbf{u}}_1^* \rangle \geq 0. \quad (184)$$

Since we have already bounded $\|\mathbf{M}^\tau - \mathbf{M}^{*\tau}\|$ in Lemma B.1, we can decompose

$$\mathbf{M}^{\tau,(m)} - \mathbf{M}^{*\tau} = \mathbf{M}^{\tau,(m)} - \mathbf{M}^\tau + \mathbf{M}^\tau - \mathbf{M}^{*\tau} = \mathbf{M}^{\tau,(m)} - \widehat{\mathbf{M}}^{\tau,(m)} + \widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^\tau + \mathbf{M}^\tau - \mathbf{M}^{*\tau}.$$

With these definitions in place, Lemma B.2 below provides the desired bounds.

Lemma B.2. *Instate the assumptions of Lemma 5.14. With probability at least $1 - O(d^{-10})$, the following holds simultaneously for all $1 \leq m \leq d$:*

$$\|\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)}\| \lesssim \frac{\mu r \lambda_{\max}^* \sqrt{\log d}}{d \sqrt{p}} + \sigma \sqrt{\frac{rd \log d}{p}}, \quad (185)$$

$$\|\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^\tau\| \lesssim \mathcal{E}_{100} \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^*, \quad (186)$$

where \mathcal{E}_{100} is defined in (65). As a result, one has

$$\|(\mathbf{M}^{\tau,(m)} - \mathbf{M}^{*\tau}) \bar{\mathbf{u}}_1^*\|_2 \leq \|\mathbf{M}^\tau - \mathbf{M}^{\tau,(m)}\| \lesssim \frac{\mu r \lambda_{\max}^* \sqrt{\log d}}{d \sqrt{p}} + \sigma \sqrt{\frac{rd \log d}{p}}. \quad (187)$$

Proof. See Appendix B.5. ■

We then can further combine (161) and (162) to deduce that

$$\begin{aligned} \|\mathbf{M}^{\tau,(m)} - \mathbf{M}^{*\tau}\| &\leq \|\mathbf{M}^\tau - \mathbf{M}^{\tau,(m)}\| + \|\mathbf{M}^\tau - \mathbf{M}^{*\tau}\| \\ &\lesssim \underbrace{\frac{\sqrt{\mu r} \lambda_{\max}^* \log^3 d}{d^{3/2} p} + \frac{\mu r \lambda_{\max}^* \log^{5/2} d}{d \sqrt{p}} + \frac{\sigma \log^{7/2} d}{p} + \sigma \sqrt{\frac{rd \log^5 d}{p}} + \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^*}_{=\mathcal{E}_{\text{op}} \cdot \lambda_{\min}^*}, \end{aligned} \quad (188)$$

and

$$\|(\mathbf{M}^{\tau,(m)} - \mathbf{M}^{*\tau}) \bar{\mathbf{u}}_1^*\|_2 \lesssim \underbrace{\frac{\mu r \lambda_{\max}^* \log d}{d \sqrt{p}} + \sigma \sqrt{\frac{rd \log^2 d}{p}} + \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^*}_{=\mathcal{E}_{\text{proj}} \cdot \lambda_{\min}^*}. \quad (189)$$

In particular, we have $\|\mathbf{M}^{\tau,(m)} - \mathbf{M}^{*\tau}\| \ll \lambda_{\min}^*$ under our conditions, and it follows from (164) that

$$\sigma_1(\mathbf{M}^{*\tau}) - \sigma_2(\mathbf{M}^{*\tau}) - \|\mathbf{M}^{\tau,(m)} - \mathbf{M}^{*\tau}\| \gtrsim \lambda_{\min}^*.$$

Invoke the bound (182) to obtain

$$\|\bar{\mathbf{u}}^{\tau,(m)} - \bar{\mathbf{u}}_1^*\|_2 \lesssim \frac{1}{\lambda_{\min}^*} \|(\mathbf{M}^{\tau,(m)} - \mathbf{M}^{*\tau}) \bar{\mathbf{u}}_1^*\|_2 \lesssim \underbrace{\frac{\mu r \log d}{d \sqrt{p}} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{rd \log^2 d}{p}} + \sqrt{\frac{\mu r \log d}{d}}}_{=\mathcal{E}_{\text{proj}}}. \quad (190)$$

To finish up, combine this with (181) and the spectral condition to arrive at

$$\alpha_3 \lesssim \mathcal{E}_{\text{proj}} \sqrt{\frac{\mu r \log d}{d}} \lesssim \mathcal{E}_{\text{op}} \sqrt{\frac{\mu r \log d}{d}}, \quad (191)$$

which results from the fact that $\mathcal{E}_{\text{proj}} \leq \mathcal{E}_{\text{op}}$ (cf. (81)).

B.4.2 Controlling α_2

We then turn to α_2 (cf. (179)). Recall the definition of $\mathbf{M}^{\tau,(m)}$ in (70c). It is straightforward to verify that

$$\begin{aligned} (\mathbf{M}^{\tau,(m)} - \mathbf{M}^{*\tau})_{m,:} &= (p^{-1} \mathbf{T}^{\tau,(m)} \times_3 \boldsymbol{\theta}^{\tau,(m)} - \mathbf{T}^* \times_3 \boldsymbol{\theta}^\tau)_{m,:} + (\mathbf{T}^* \times_3 \boldsymbol{\theta}^\tau - \mathbf{M}^{*\tau})_{m,:} \\ &= (p^{-1} \mathbf{T}^{\tau,(m)} \times_3 \boldsymbol{\theta}^{\tau,(m)} - \mathbf{T}^* \times_3 \boldsymbol{\theta}^\tau)_{m,:} + \mathbf{C}_{m,:}^\tau \\ &= \mathbf{T}_{m,:}^* \times_3 (\boldsymbol{\theta}^{\tau,(m)} - \boldsymbol{\theta}^\tau) + \mathbf{C}_{m,:}^\tau, \end{aligned}$$

where \mathbf{C}^τ is defined in (156).

From the incoherence conditions, we can upper bound

$$\begin{aligned} \|\mathbf{C}_{m,:}^\tau\|_2 &= \left\| \sum_{s:s \neq 1} \gamma_s^{*\tau} \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle ((\bar{\mathbf{u}}_1^*)_m \bar{\mathbf{u}}_s^{*\top} + (\bar{\mathbf{u}}_s^*)_m \bar{\mathbf{u}}_1^{*\top}) - \sum_{s:s \neq 1} \gamma_s^{*\tau} \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle^2 (\bar{\mathbf{u}}_1^*)_m \bar{\mathbf{u}}_1^{*\top} \right\|_2 \\ &\leq \left\| \sum_{s:s \neq 1} \gamma_s^{*\tau} \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle (\bar{\mathbf{u}}_1^*)_m \bar{\mathbf{u}}_s^{*\top} \right\|_2 + \left| \sum_{s:s \neq 1} \gamma_s^{*\tau} \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle (\bar{\mathbf{u}}_s^*)_m \right| + \left| \sum_{s:s \neq 1} \gamma_s^{*\tau} \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle^2 (\bar{\mathbf{u}}_1^*)_m \right| \\ &\lesssim \|\boldsymbol{\gamma}^{*\tau}\|_2 \left(\max_{s:s \neq 1} |\langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle| \|\bar{\mathbf{u}}_1^*\|_\infty \|\bar{\mathbf{U}}^*\| + \max_{s:s \neq 1} |\langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle| \|\bar{\mathbf{U}}^*\|_{2,\infty} + \max_{s:s \neq 1} |\langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle|^2 \|\bar{\mathbf{u}}_1^*\|_\infty \sqrt{r} \right) \\ &\lesssim \|\boldsymbol{\gamma}^{*\tau}\|_2 \left(\sqrt{\frac{\mu}{d}} \sqrt{\frac{\mu}{d}} + \sqrt{\frac{\mu}{d}} \sqrt{\frac{\mu r}{d}} + \sqrt{\frac{\mu}{d}} \sqrt{\frac{\mu r}{d}} \right) \lesssim \sqrt{\frac{\mu r}{d}} \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^*, \end{aligned} \quad (192)$$

where we use the fact that $\|\boldsymbol{\gamma}^{*\tau}\|_2 \lesssim \sqrt{r \log d} \lambda_{\max}^*$ from Lemma 5.12 and $\|\bar{\mathbf{U}}^*\| \lesssim 1$ from Lemma D.1.

In addition, we can express

$$\mathbf{T}_{m,:}^* \times_3 (\boldsymbol{\theta}^{\tau,(m)} - \boldsymbol{\theta}^\tau) = \sum_{s \in [r]} \lambda_s^* (\bar{\mathbf{u}}_s^*)_m \langle \bar{\mathbf{u}}_s^*, \boldsymbol{\theta}^{\tau,(m)} - \boldsymbol{\theta}^\tau \rangle \bar{\mathbf{u}}_s^*.$$

By construction, we know that

$$\boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{\tau,(m)} = (\mathbf{U}\mathbf{U}^\top - \mathbf{U}^{(m)}\mathbf{U}^{(m)\top}) \mathbf{g}^\tau$$

is a zero-mean Gaussian random vector conditional on $\mathcal{P}_\Omega(\mathbf{E})$. Using standard results on Gaussian random vectors, one has: with probability at least $1 - O(d^{-11})$, for each $s \in [r]$ and $m \in [d]$,

$$\begin{aligned} \left| \langle \bar{\mathbf{u}}_s^*, \boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{\tau,(m)} \rangle \right| &= \left| \langle \mathbf{g}^\tau, (\mathbf{U}\mathbf{U}^\top - \mathbf{U}^{(m)}\mathbf{U}^{(m)\top}) \bar{\mathbf{u}}_s^* \rangle \right| \lesssim \|(\mathbf{U}\mathbf{U}^\top - \mathbf{U}^{(m)}\mathbf{U}^{(m)\top}) \bar{\mathbf{u}}_s^*\|_2 \sqrt{\log d} \\ &\leq \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^{(m)}\mathbf{U}^{(m)\top}\| \sqrt{\log d} \end{aligned} \quad (193)$$

and

$$\|\boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{\tau,(m)}\|_2 \lesssim \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^{(m)}\mathbf{U}^{(m)\top}\|_{\text{F}} \sqrt{\log d}. \quad (194)$$

Therefore, we have

$$\begin{aligned} \|\mathbf{T}_{m,:}^* \times_3 (\boldsymbol{\theta}^{\tau,(m)} - \boldsymbol{\theta}^\tau)\|_2 &\leq \max_{s \in [r]} |\lambda_s^* \langle \bar{\mathbf{u}}_s^*, \boldsymbol{\theta}^{\tau,(m)} - \boldsymbol{\theta}^\tau \rangle| \|\bar{\mathbf{U}}^*\|_{2,\infty} \|\bar{\mathbf{U}}^*\| \\ &\lesssim \lambda_{\max}^* \|\bar{\mathbf{U}}^*\|_{2,\infty} \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^{(m)}\mathbf{U}^{(m)\top}\| \sqrt{\log d} \\ &\lesssim \mathcal{E}_{\text{loo}} \sqrt{\frac{\mu r}{d}} \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^* \ll \sqrt{\frac{\mu r}{d}} \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^* \end{aligned}$$

where we have used Lemma 5.8 as well as the conditions that $\mathcal{E}_{\text{loo}} \ll 1$, $\|\bar{\mathbf{U}}^*\| \lesssim 1$ and $\|\bar{\mathbf{U}}^*\|_{2,\infty} \lesssim \sqrt{\mu r/d}$.

Putting the above bounds together, we arrive at

$$\|(\mathbf{M}^{\tau,(m)} - \mathbf{M}^{*\tau})_{m,:}\|_2 \leq \|\mathbf{T}_{m,:}^* \times_3 (\boldsymbol{\theta}^{\tau,(m)} - \boldsymbol{\theta}^\tau)\|_2 + \|\mathbf{C}_{m,:}^\tau\|_2 \lesssim \sqrt{\frac{\mu r}{d}} \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^*. \quad (195)$$

We therefore conclude that

$$\alpha_2 \lesssim \sqrt{\frac{\mu r}{d}} \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^* \lesssim \mathcal{E}_{\text{op}} \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^*$$

where we remind the reader of the definition of \mathcal{E}_{op} in (81).

B.4.3 Controlling α_1

The remaining quantity to control is α_1 (see (179)). Invoke Weyl's inequality to show that

$$|\gamma_1^{*\tau} - \gamma_\tau^{(m)}| \leq \|\mathbf{M}^{\tau,(m)} - \mathbf{M}^{*\tau}\| \lesssim \mathcal{E}_{\text{op}} \lambda_{\min}^* \ll \gamma_1^{*\tau},$$

where the last inequality arises from (188) and Lemma 5.12. Under our sample size, rank and noise conditions, we have

$$\frac{1}{2} \gamma_1^{*\tau} \leq \gamma_1^{*\tau} - |\gamma_1^{*\tau} - \gamma_\tau^{(m)}| \leq \gamma_\tau^{(m)} \leq |\gamma_1^{*\tau} - \gamma_\tau^{(m)}| + \gamma_1^{*\tau} \leq 2\gamma_1^{*\tau}.$$

This indicates that

$$\frac{|\gamma_1^{*\tau} - \gamma_\tau^{(m)}|}{\gamma_1^{*\tau} \gamma_\tau^{(m)}} \lesssim \frac{1}{(\gamma_1^{*\tau})^2} \|\mathbf{M}^{\tau,(m)} - \mathbf{M}^{*\tau}\| \lesssim \frac{1}{\lambda_{\min}^*} \mathcal{E}_{\text{op}}.$$

Moreover, we learn from (195) and (181) that

$$\|\mathbf{M}_{m,:}^{\tau,(m)}\|_2 \leq \|(\mathbf{M}^{\tau,(m)} - \mathbf{M}^{*\tau})_{m,:}\|_2 + \|\mathbf{M}_{m,:}^{*\tau}\|_2$$

$$\lesssim \frac{\mu r \sqrt{\log d}}{d} \lambda_{\max}^* + \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^* \asymp \sqrt{\frac{\mu r \log d}{d}} \lambda_{\min}^*,$$

where the last step follows from the fact that $\mu r \leq d$ and $\kappa \asymp 1$. Hence, we reach the conclusion that

$$\alpha_1 \lesssim \frac{|\gamma_1^{*\tau} - \gamma_\tau^{(m)}|}{\gamma_1^{*\tau} \gamma_\tau^{(m)}} \|\mathbf{M}_{m,:}^{\tau,(m)}\|_2 \lesssim \mathcal{E}_{\text{op}} \sqrt{\frac{\mu r \log d}{d}}.$$

B.4.4 Combining α_1 , α_2 and α_3

Putting together all of the preceding bounds on α_1 , α_2 and α_3 immediately establishes the lemma.

B.5 Proof of Lemma B.2

First of all, if the claims (185) and (186) can be established, then putting them together yields

$$\|\mathbf{M}^\tau - \mathbf{M}^{\tau,(m)}\| \leq \|\mathbf{M}^\tau - \widehat{\mathbf{M}}^{\tau,(m)}\| + \|\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)}\| \lesssim \frac{\mu r \lambda_{\max}^* \sqrt{\log d}}{d \sqrt{p}} + \sigma \sqrt{\frac{rd \log d}{p}}, \quad (196)$$

where we recall the definition of \mathcal{E}_{loo} in (65) and use the sample size, noise and rank conditions. The rest of the proof is thus dedicated to establishing (185) and (186). In what follows, we shall assume $\{E_{i,j,k}\}_{i,j,k \in [d]}$ (resp. $\{\chi_{i,j,k}\}_{i,j,k \in [d]}$) are independent random variables to simplify presentation.

B.5.1 Proximity of $\mathbf{M}^{\tau,(m)}$ and $\widehat{\mathbf{M}}^{\tau,(m)}$

Recall the definition of $\mathbf{M}^{\tau,(m)} = p^{-1} \mathbf{T}^{(m)} \times_3 \boldsymbol{\theta}^{\tau,(m)}$ in (70c). Comparing this with the definition of $\widehat{\mathbf{M}}^{\tau,(m)}$ in (183), we see that

$$(\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)})_{i,j} = \theta_m^{\tau,(m)} (T_{i,j,m}^* (p^{-1} \chi_{i,j,m} - 1) + p^{-1} E_{i,j,m} \chi_{i,j,m}), \quad i \neq m, j \neq m, \quad (197a)$$

$$(\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)})_{i,m} = \theta_k^{\tau,(m)} \sum_{k \in [d]} (T_{i,m,k}^* (p^{-1} \chi_{i,m,k} - 1) + p^{-1} E_{i,m,k} \chi_{i,m,k}), \quad i \neq m, \quad (197b)$$

$$(\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)})_{m,j} = \theta_k^{\tau,(m)} \sum_{k \in [d]} (T_{m,j,k}^* (p^{-1} \chi_{m,j,k} - 1) + p^{-1} E_{m,j,k} \chi_{m,j,k}). \quad (197c)$$

Note that $\boldsymbol{\theta}^{\tau,(m)} \sim \mathcal{N}(\mathbf{0}, \mathbf{U}^{(m)} \mathbf{U}^{(m)\top})$ conditional on $\mathcal{P}_\Omega(\mathbf{E})$. Standard Gaussian concentration inequalities reveal that with probability exceeding $1 - O(d^{-10})$,

$$\|\boldsymbol{\theta}^{\tau,(m)}\|_2 \lesssim \sqrt{r \log d}. \quad (198)$$

From Lemmas 5.7-5.8 and the fact that $\max\{\mathcal{E}_{\text{se}}, \mathcal{E}_{\text{loo}}\} \ll 1$, we have

$$\begin{aligned} \max_{i \in [d]} \text{Var}(\theta_i^\tau) &= \|\mathbf{U}\|_{2,\infty}^2 \lesssim \frac{\mu r}{d} + \|\mathbf{U}_{\text{orth}}^*\|_{2,\infty}^2 \asymp \frac{\mu r}{d}, \\ \max_{i \in [d]} \text{Var}(\theta_i^{\tau,(m)}) &= \|\mathbf{U}^{(m)}\|_{2,\infty}^2 \lesssim \|\mathbf{U}\|_{2,\infty}^2 \lesssim \frac{\mu r}{d}, \quad 1 \leq m \leq d. \end{aligned}$$

As a consequence, standard concentration results assert that with probability $1 - O(d^{-10})$,

$$\|\boldsymbol{\theta}^\tau\|_\infty \leq \sqrt{\max_{i \in [d]} \text{Var}(\theta_i^\tau) \log d} \lesssim \sqrt{\frac{\mu r \log d}{d}}; \quad (199)$$

$$\|\boldsymbol{\theta}^{\tau,(m)}\|_\infty \leq \sqrt{\max_{i \in [d]} \text{Var}(\theta_i^{\tau,(m)}) \log d} \lesssim \sqrt{\frac{\mu r \log d}{d}}, \quad 1 \leq m \leq d. \quad (200)$$

- Regarding the m -th row of $\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)}$, apply Lemma D.9 to show that with probability $1 - O(d^{-11})$,

$$\sum_{j \in [d]} \left(\sum_{k \in [d]} T_{m,j,k}^* \theta_k^{\tau,(m)} (p^{-1} \chi_{m,j,k} - 1) \right)^2 \lesssim \frac{\mu r \lambda_{\max}^{*2}}{dp} \|\boldsymbol{\theta}^{\tau,(m)}\|_{\infty}^2 \lesssim \frac{\mu^2 r^2 \lambda_{\max}^{*2} \log d}{d^2 p},$$

where the last inequality comes from (200). In addition, Lemma D.10 indicates that with probability exceeding $1 - O(d^{-11})$,

$$\begin{aligned} \sum_{j \in [d]} \left(\sum_{k \in [d]} \theta_k^{\tau,(m)} E_{m,j,k} \chi_{m,j,k} \right)^2 &\lesssim \sigma^2 dp \|\boldsymbol{\theta}^{\tau,(m)}\|_2^2 + \sigma^2 \|\boldsymbol{\theta}^{\tau,(m)}\|_{\infty}^2 \log^5 d \\ &\lesssim \sigma^2 r dp \log d + \frac{\sigma^2 \mu r \log^6 d}{d} \asymp \sigma^2 r dp \log d, \end{aligned}$$

where the second line comes from (200) and (198), and the last inequality holds as long as $p \gg \mu d^{-2} \log^5 d$. These together with (197c) allow us to obtain

$$\left\| \left(\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)} \right)_{m,:} \right\|_2^2 \lesssim \frac{\mu^2 r^2 \lambda_{\max}^{*2} \log d}{d^2 p} + \frac{\sigma^2 r d \log d}{p}. \quad (201)$$

Clearly, this bound is also valid for $\sum_{i:i \neq m} \left\{ \left(\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)} \right)_{i,m} \right\}^2$, namely,

$$\sum_{i:i \neq m} \left\{ \left(\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)} \right)_{i,m} \right\}^2 \lesssim \frac{\mu^2 r^2 \lambda_{\max}^{*2} \log d}{d^2 p} + \frac{\sigma^2 r d \log d}{p}. \quad (202)$$

- When it comes to the remaining entries of $\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)}$, by the fact that the spectral norm of a submatrix is always less than or equal to that of the whole matrix, applying the matrix Bernstein inequality gives that with probability $1 - O(d^{-11})$,

$$\begin{aligned} \left\| \left[\left(\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)} \right)_{i,j} \right]_{i,j \neq m} \right\| &\lesssim \|\boldsymbol{\theta}^{\tau,(m)}\|_{\infty} \left(\frac{\log d}{p} \|\mathbf{A}^*\|_{\infty} + \sqrt{\frac{\log d}{p}} \|\mathbf{A}^{*\top}\|_{2,\infty} + \frac{\sigma \log^2 d}{p} + \sigma \sqrt{\frac{d \log d}{p}} \right) \\ &\lesssim \sqrt{\frac{\mu r \log d}{d}} \left(\frac{\sqrt{\mu r} \lambda_{\max}^* \log d}{d^{3/2} p} + \frac{\mu \sqrt{r} \lambda_{\max}^* \sqrt{\log d}}{d \sqrt{p}} + \frac{\sigma \log^2 d}{p} + \sigma \sqrt{\frac{d \log d}{p}} \right) \\ &\lesssim \frac{\mu r \lambda_{\max}^* \sqrt{\log d}}{d \sqrt{p}} + \sigma \sqrt{\frac{rd \log d}{p}}, \end{aligned}$$

as long as our sample size and rank condition holds.

- Putting the preceding bounds together yields

$$\left\| \widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)} \right\| \lesssim \frac{\mu r \lambda_{\max}^* \sqrt{\log d}}{d \sqrt{p}} + \sigma \sqrt{\frac{rd \log d}{p}}. \quad (203)$$

B.5.2 Proximity of \mathbf{M}^{τ} and $\widehat{\mathbf{M}}^{\tau,(m)}$

Recall the definitions of \mathbf{M}^{τ} and $\widehat{\mathbf{M}}^{\tau,(m)}$ in (279b) and (183), respectively. From the definition of the operator norm and the triangle inequality, we have

$$\left\| \mathbf{M}^{\tau} - \widehat{\mathbf{M}}^{\tau,(m)} \right\| \leq \underbrace{\left\| \mathbf{T}^* \times_3 (\boldsymbol{\theta}^{\tau} - \boldsymbol{\theta}^{\tau,(m)}) \right\|}_{=: \alpha_1} + \underbrace{\left\| (p^{-1} \mathbf{T} - \mathbf{T}^*) \times_3 (\boldsymbol{\theta}^{\tau} - \boldsymbol{\theta}^{\tau,(m)}) \right\|}_{=: \alpha_2}. \quad (204)$$

- To control α_1 , we can express

$$\mathbf{T}^* \times_3 (\boldsymbol{\theta}^{\tau} - \boldsymbol{\theta}^{\tau,(m)}) = \sum_{s \in [r]} \lambda_s^* \langle \mathbf{u}_s^*, \boldsymbol{\theta}^{\tau} - \boldsymbol{\theta}^{\tau,(m)} \rangle \mathbf{u}_s^* \mathbf{u}_s^{*\top}. \quad (205)$$

As shown in (193), with probability at least $1 - O(d^{-12})$,

$$|\langle \bar{\mathbf{u}}_s^*, \boldsymbol{\theta}^{\tau, (m)} - \boldsymbol{\theta}^\tau \rangle| \lesssim \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^{(m)}\mathbf{U}^{(m)\top}\| \sqrt{\log d}.$$

Consequently, we know from Lemma 5.8 that with probability at least $1 - O(d^{-10})$,

$$\begin{aligned} \|\mathbf{T}^* \times_3 (\boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{\tau, (m)})\| &\leq \max_{s \in [r]} |\lambda_s^* \langle \bar{\mathbf{u}}_s^*, \boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{\tau, (m)} \rangle| \|\bar{\mathbf{U}}^*\|^2 \lesssim \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^{(m)}\mathbf{U}^{(m)\top}\| \sqrt{\log d} \lambda_{\max}^* \\ &\lesssim \mathcal{E}_{\text{loo}} \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^*, \end{aligned} \quad (206)$$

where we use the fact that $\|\bar{\mathbf{U}}^*\| \lesssim 1$ if $r\sqrt{\mu/d} \ll 1$.

- When it comes to α_2 , combining (171) and (194) with our sample size, noise and rank conditions, one has

$$\begin{aligned} \|(p^{-1}\mathbf{T} - \mathbf{T}^*) \times_3 (\boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{\tau, (m)})\| &\leq \|p^{-1}\mathbf{T} - \mathbf{T}^*\| \|\boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{\tau, (m)}\|_2 \ll \lambda_{\max}^* \|\boldsymbol{\theta}^\tau - \boldsymbol{\theta}^{\tau, (m)}\|_2 \\ &\lesssim \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^{(m)}\mathbf{U}^{(m)\top}\|_{\text{F}} \sqrt{\log d} \lambda_{\max}^* \\ &\lesssim \mathcal{E}_{\text{loo}} \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^*. \end{aligned} \quad (207)$$

- Combining (204), (206) and (207), we conclude that

$$\|\mathbf{M}^\tau - \widehat{\mathbf{M}}^{\tau, (m)}\| \lesssim \mathcal{E}_{\text{loo}} \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^*.$$

B.6 Proof of Lemma 5.15

- We start with the first claim regarding $\|\bar{\mathbf{u}}^\tau - \boldsymbol{\nu}^{\tau, (m)}\|_2$, or equivalently, $\|\bar{\mathbf{u}}^\tau - \bar{\mathbf{u}}^{\tau, (m)}\|_2$ (as argued in the proof of Lemma 5.14). By the triangle inequality, we can upper bound the following two terms separately:

$$\|\bar{\mathbf{u}}^\tau - \bar{\mathbf{u}}^{\tau, (m)}\|_2 \leq \underbrace{\|\bar{\mathbf{u}}^\tau - \hat{\mathbf{u}}^{\tau, (m)}\|_2}_{=: \beta_1} + \underbrace{\|\hat{\mathbf{u}}^{\tau, (m)} - \bar{\mathbf{u}}^{\tau, (m)}\|_2}_{=: \beta_2}. \quad (208)$$

Here, we remind the reader that $\hat{\mathbf{u}}^{\tau, (m)}$ is the top left singular vector of $\widehat{\mathbf{M}}^{\tau, (m)}$ (see (183)) obeying $\langle \hat{\mathbf{u}}^{\tau, (m)}, \bar{\mathbf{u}}_1^* \rangle \geq 0$.

- The first term β_1 shall be bounded via Wedin's theorem. From (163) and (164), we have

$$\sigma_1(\mathbf{M}^\tau) - \sigma_2(\mathbf{M}^\tau) \geq \sigma_1(\mathbf{M}^{*\tau}) - \sigma_2(\mathbf{M}^{*\tau}) - 2\|\mathbf{M}^\tau - \mathbf{M}^{*\tau}\| \gtrsim \lambda_{\min}^*. \quad (209)$$

Combined with Lemma B.2, one has

$$\sigma_1(\mathbf{M}^\tau) - \sigma_2(\mathbf{M}^\tau) - \|\mathbf{M}^\tau - \widehat{\mathbf{M}}^{\tau, (m)}\| \gtrsim \lambda_{\min}^*.$$

Note that we have already shown in the proof of Lemma B.2 and Lemma B.4 that $\|\bar{\mathbf{u}}^\tau - \bar{\mathbf{u}}_1^*\|_2 = o(1)$ and $\|\hat{\mathbf{u}}^{\tau, (m)} - \bar{\mathbf{u}}_1^*\|_2 = o(1)$, which implies that \mathbf{u}^τ and $\hat{\mathbf{u}}^{\tau, (m)}$ are positively correlated. Thus, one can invoke Wedin's theorem and use the bound (186) to reach

$$\begin{aligned} \|\bar{\mathbf{u}}^\tau - \hat{\mathbf{u}}^{\tau, (m)}\|_2 &\leq \frac{\|\mathbf{M}^\tau - \widehat{\mathbf{M}}^{\tau, (m)}\|}{\sigma_1(\mathbf{M}^\tau) - \sigma_2(\mathbf{M}^\tau) - \|\mathbf{M}^\tau - \widehat{\mathbf{M}}^{\tau, (m)}\|} \lesssim \frac{1}{\lambda_{\min}^*} \|\mathbf{M}^\tau - \widehat{\mathbf{M}}^{\tau, (m)}\| \\ &\lesssim \mathcal{E}_{\text{loo}} \sqrt{\frac{\mu r \log d}{d}} \leq \mathcal{E}_{\text{loo}} \sqrt{r \log d} \max\{\sqrt{\mu/d}, \|\bar{\mathbf{u}}^\tau\|_\infty\}. \end{aligned} \quad (210)$$

In addition to this bound on β_1 , we also make note of the following simple bound

$$\|\hat{\mathbf{u}}^{\tau, (m)}\|_\infty \leq \|\bar{\mathbf{u}}^\tau - \hat{\mathbf{u}}^{\tau, (m)}\|_\infty + \|\bar{\mathbf{u}}^\tau\|_\infty \leq \|\bar{\mathbf{u}}^\tau - \hat{\mathbf{u}}^{\tau, (m)}\|_2 + \|\bar{\mathbf{u}}^\tau\|_\infty \lesssim \max\{\sqrt{\mu/d}, \|\bar{\mathbf{u}}^\tau\|_\infty\}, \quad (211)$$

where the last inequality follows from our sample size, noise and rank condition that $\mathcal{E}_{\text{loo}} \sqrt{r \log d} \ll 1$.

– The second term β_2 is also controlled via Wedin’s theorem:

$$\|\widehat{\mathbf{u}}^{\tau,(m)} - \overline{\mathbf{u}}^{\tau,(m)}\|_2 \leq \frac{\|(\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)})\overline{\mathbf{u}}^{\tau,(m)}\|_2}{\sigma_1(\mathbf{M}^{\tau,(m)}) - \sigma_2(\mathbf{M}^{\tau,(m)}) - \|\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)}\|}. \quad (212)$$

The denominator term is easy to handle. With (196) and (209) in mind, we can apply Weyl’s inequality to obtain

$$\sigma_1(\mathbf{M}^{\tau,(m)}) - \sigma_2(\mathbf{M}^{\tau,(m)}) \geq \sigma_1(\mathbf{M}^\tau) - \sigma_2(\mathbf{M}^\tau) - 2\|\mathbf{M}^\tau - \mathbf{M}^{\tau,(m)}\| \gtrsim \lambda_{\min}^*. \quad (213)$$

From Lemma B.2, one has $\|\mathbf{M}^\tau - \mathbf{M}^{\tau,(m)}\| \ll \lambda_{\min}^*$. Therefore, we know that

$$\sigma_1(\mathbf{M}^{\tau,(m)}) - \sigma_2(\mathbf{M}^{\tau,(m)}) - \|\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)}\| \gtrsim \lambda_{\min}^*.$$

In addition, Lemma B.3 below develops an upper bound on the numerator term:

Lemma B.3. *Instate the assumptions of Lemma 5.15. With probability at least $1 - O(d^{-10})$, one has*

$$\|(\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)})\overline{\mathbf{u}}^{\tau,(m)}\|_2 \lesssim \left\{ \frac{\mu r \lambda_{\max}^* \log d}{d\sqrt{p}} + \sigma \sqrt{\frac{rd \log^2 d}{p}} \right\} \|\overline{\mathbf{u}}^{\tau,(m)}\|_\infty. \quad (214)$$

Proof. See Appendix B.7. ■

– Substitution of the above bounds into (212) yields

$$\|\widehat{\mathbf{u}}^{\tau,(m)} - \overline{\mathbf{u}}^{\tau,(m)}\|_2 \lesssim \left\{ \frac{\mu r \log d}{d\sqrt{p}} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{rd \log^2 d}{p}} \right\} \|\overline{\mathbf{u}}^{\tau,(m)}\|_\infty \ll \|\overline{\mathbf{u}}^{\tau,(m)}\|_\infty,$$

where the last step holds as long as $p \gg \mu^2 r^2 d^{-2} \log^2 d$ and $\sigma/\lambda_{\min}^* \ll \sqrt{p/(rd \log^2 d)}$. In addition, from (211), we observe that

$$\begin{aligned} \|\overline{\mathbf{u}}^{\tau,(m)}\|_\infty &\leq \|\widehat{\mathbf{u}}^{\tau,(m)} - \overline{\mathbf{u}}^{\tau,(m)}\|_\infty + \|\widehat{\mathbf{u}}^{\tau,(m)}\|_\infty \leq \|\widehat{\mathbf{u}}^{\tau,(m)} - \overline{\mathbf{u}}^{\tau,(m)}\|_2 + \|\widehat{\mathbf{u}}^{\tau,(m)}\|_\infty \\ &\lesssim o(1) \|\overline{\mathbf{u}}^{\tau,(m)}\|_\infty + \|\widehat{\mathbf{u}}^{\tau,(m)}\|_\infty \leq o(1) \|\overline{\mathbf{u}}^{\tau,(m)}\|_\infty + \max\{\sqrt{\mu/d}, \|\overline{\mathbf{u}}^\tau\|_\infty\}, \end{aligned}$$

from which we can deduce that

$$\|\overline{\mathbf{u}}^{\tau,(m)}\|_\infty \lesssim \max\{\sqrt{\mu/d}, \|\overline{\mathbf{u}}^\tau\|_\infty\}.$$

As a consequence, one immediately obtains

$$\|\widehat{\mathbf{u}}^{\tau,(m)} - \overline{\mathbf{u}}^{\tau,(m)}\|_2 \lesssim \left\{ \frac{\mu r \log d}{d\sqrt{p}} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{rd \log^2 d}{p}} \right\} \max\{\sqrt{\mu/d}, \|\overline{\mathbf{u}}^\tau\|_\infty\}. \quad (215)$$

– Combining (208), (210) and (215) and the definition of \mathcal{E}_{loo} , we arrive at

$$\|\overline{\mathbf{u}}^\tau - \overline{\mathbf{u}}^{\tau,(m)}\|_2 \lesssim \mathcal{E}_{\text{loo}} \sqrt{r \log d} \max\{\sqrt{\mu/d}, \|\overline{\mathbf{u}}^\tau\|_\infty\}. \quad (216)$$

Comparing this bound with the first claim of the lemma, we see that the claim can be established as long as we can show that

$$\|\overline{\mathbf{u}}^\tau\|_\infty \lesssim \sqrt{\mu/d}. \quad (217)$$

To justify this bound (217), we make use of Lemma 5.14 to derive that

$$|(\overline{\mathbf{u}}^\tau - \overline{\mathbf{u}}_1^*)|_m \leq |(\overline{\mathbf{u}}^\tau - \overline{\mathbf{u}}^{\tau,(m)})|_m + |(\overline{\mathbf{u}}^{\tau,(m)} - \overline{\mathbf{u}}_1^*)|_m$$

$$\begin{aligned} &\leq \|\bar{\mathbf{u}}^\tau - \bar{\mathbf{u}}^{\tau,(m)}\|_2 + \left| (\bar{\mathbf{u}}^{\tau,(m)} - \bar{\mathbf{u}}_1^*)^m \right| \\ &\lesssim (\mathcal{E}_{\text{loo}} + \mathcal{E}_{\text{op}}) \sqrt{r \log d} \max \{ \sqrt{\mu/d}, \|\bar{\mathbf{u}}^\tau\|_\infty \} \end{aligned}$$

for each $m \in [d]$. Maximizing over $m \in [d]$ gives that

$$\|\bar{\mathbf{u}}^\tau - \bar{\mathbf{u}}_1^*\|_\infty \lesssim (\mathcal{E}_{\text{loo}} + \mathcal{E}_{\text{op}}) \sqrt{r \log d} \max \{ \sqrt{\mu/d}, \|\bar{\mathbf{u}}^\tau\|_\infty \} \quad (218)$$

$$\ll \max \{ \sqrt{\mu/d}, \|\bar{\mathbf{u}}^\tau\|_\infty \}, \quad (219)$$

where we use the condition that $(\mathcal{E}_{\text{loo}} + \mathcal{E}_{\text{op}}) \sqrt{r \log d} \ll 1$. Apply the triangle inequality to yield

$$\|\bar{\mathbf{u}}^\tau\|_\infty \leq \|\bar{\mathbf{u}}^\tau - \bar{\mathbf{u}}_1^*\|_\infty + \|\bar{\mathbf{u}}_1^*\|_\infty \leq o(1) \|\bar{\mathbf{u}}^\tau\|_\infty + \sqrt{\mu/d}.$$

These allow us to establish the claim (217), which in turn finishes the proof for the first claim of this lemma.

- The second claim (83) of this lemma follows immediately from (217) and (218).
- It remains to prove the last claim (84). Recall the definition of λ_τ and $\lambda_\tau^{(m)}$ in (71). We can decompose

$$\begin{aligned} \langle p^{-1} \mathbf{T}^{(m)}, (\bar{\mathbf{u}}^{\tau,(m)})^{\otimes 3} \rangle - \langle p^{-1} \mathbf{T}, (\bar{\mathbf{u}}^\tau)^{\otimes 3} \rangle &= \underbrace{\langle p^{-1} \mathbf{T}^{(m)} - p^{-1} \mathbf{T}, (\bar{\mathbf{u}}^{\tau,(m)})^{\otimes 3} \rangle}_{=: \beta_1} \\ &\quad + \underbrace{\langle p^{-1} \mathbf{T}, (\bar{\mathbf{u}}^{\tau,(m)})^{\otimes 3} - (\bar{\mathbf{u}}^\tau)^{\otimes 3} \rangle}_{=: \beta_2} \end{aligned} \quad (220)$$

In what follows, we will control β_1 and β_2 separately.

For β_1 , we note that all non-zero entries of $\mathbf{T}^{(m)} - \mathbf{T}$ are located in the m th slices, and are independent of $\bar{\mathbf{u}}^{\tau,(m)}$. This type of quantities have appeared many times and we omit the detailed proof for conciseness. By the Bernstein inequality, one can show that with probability at least $1 - O(d^{-10})$,

$$\left| \langle p^{-1} \mathbf{T}^{(m)} - p^{-1} \mathbf{T}, (\bar{\mathbf{u}}^{\tau,(m)})^{\otimes 3} \rangle \right| \lesssim \sqrt{\frac{\mu r \log d}{d^2 p}} \sqrt{\frac{\mu}{d}} \lambda_{\max}^*. \quad (221)$$

Next, we turn to β_2 . From our sample size and noise condition, Lemma D.1 and Corollary D.3 demonstrates that with probability at least $1 - O(d^{-10})$,

$$\|p^{-1} \mathbf{T}\| \leq \|p^{-1} \mathbf{T} - \mathbf{T}^*\| + \|\mathbf{T}^*\| \leq \|p^{-1} \mathbf{T} - \mathbf{T}^*\| + \|\mathbf{A}^*\| \lesssim \lambda_{\max}^*,$$

where we use the fact that the tensor spectral norm is always less than or equal to that of its matricization. By the definition of the operator norm, one has

$$\begin{aligned} |\beta_2| &\leq 3 \left| \langle p^{-1} \mathbf{T}, (\bar{\mathbf{u}}^{\tau,(m)})^{\otimes 2} \otimes (\bar{\mathbf{u}}^{\tau,(m)} - \bar{\mathbf{u}}^\tau) \rangle \right| + 3 \left| \langle p^{-1} \mathbf{T}, (\bar{\mathbf{u}}^{\tau,(m)}) \otimes (\bar{\mathbf{u}}^{\tau,(m)} - \bar{\mathbf{u}}^\tau)^{\otimes 2} \rangle \right| \\ &\quad + \left| \langle p^{-1} \mathbf{T}, (\bar{\mathbf{u}}^{\tau,(m)} - \bar{\mathbf{u}}^\tau)^{\otimes 3} \rangle \right| \\ &\lesssim \|p^{-1} \mathbf{T}\| \left(\|\bar{\mathbf{u}}^{\tau,(m)} - \bar{\mathbf{u}}^\tau\|_2 + \|\bar{\mathbf{u}}^{\tau,(m)} - \bar{\mathbf{u}}^\tau\|_2^2 + \|\bar{\mathbf{u}}^{\tau,(m)} - \bar{\mathbf{u}}^\tau\|_2^3 \right) \\ &\lesssim \mathcal{E}_{\text{loo}} \sqrt{\frac{\mu r \log d}{d}} \lambda_{\max}^* \end{aligned} \quad (222)$$

where we use (216) and (217) in the last step and the fact that $\mathcal{E}_{\text{loo}} \sqrt{\mu r \log d/d} \ll 1$.

Combining (221) and (222) immediately establishes (84).

B.7 Proof of Lemma B.3

Recalling the definitions of $\widehat{\mathbf{M}}^{\tau,(m)}$ and $\mathbf{M}^{\tau,(m)}$ in (183) and (70c), respectively, we observe that $\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)}$ is independent of $\overline{\mathbf{u}}^{\tau,(m)}$ conditional on $\mathcal{P}_{\Omega_{-m}}(\mathbf{E})$ and \mathbf{g} .

- The m -th entry of $(\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)})\overline{\mathbf{u}}^{\tau,(m)}$ can be written as

$$\begin{aligned} & (\widehat{\mathbf{M}}^{\tau,(m)} - \mathbf{M}^{\tau,(m)})_{m,:} \overline{\mathbf{u}}^{\tau,(m)} = (p^{-1}(\mathbf{T} - \mathbf{T}^{(m)}) \times_3 \boldsymbol{\theta}^{\tau,(m)})_{m,:} \overline{\mathbf{u}}^{\tau,(m)} \\ & = \underbrace{\sum_{j,k \in [d]} T_{m,j,k}^* \overline{u}_j^{\tau,(m)} \theta_k^{\tau,(m)} (p^{-1} \chi_{m,j,k} - 1)}_{=: \alpha_1} + \underbrace{\sum_{j,k \in [d]} \overline{u}_j^{\tau,(m)} \theta_k^{\tau,(m)} p^{-1} E_{m,j,k} \chi_{m,j,k}}_{=: \alpha_2}. \end{aligned} \quad (223)$$

- For the first term α_1 , it is easily seen from (200) and incoherence conditions that

$$\begin{aligned} L_1 & := \max_{j,k \in [d]} \left| T_{m,j,k}^* \overline{u}_j^{\tau,(m)} \theta_k^{\tau,(m)} (p^{-1} \chi_{m,j,k} - 1) \right| \leq p^{-1} \|\mathbf{T}^*\|_{\infty} \|\overline{\mathbf{u}}^{\tau,(m)}\|_{\infty} \|\boldsymbol{\theta}^{\tau,(m)}\|_{\infty} \\ & \lesssim \frac{\mu r \lambda_{\max}^* \sqrt{\log d}}{d^2 p} \|\overline{\mathbf{u}}^{\tau,(m)}\|_{\infty}; \end{aligned}$$

and

$$\begin{aligned} V_1 & := \sum_{j,k \in [d]} T_{m,j,k}^{*2} (\overline{u}_j^{\tau,(m)})^2 (\theta_k^{\tau,(m)})^2 \mathbb{E}[(p^{-1} \chi_{m,j,k} - 1)^2] \leq p^{-1} \|\mathbf{A}^*\|_{2,\infty}^2 \|\overline{\mathbf{u}}^{\tau,(m)}\|_{\infty}^2 \|\boldsymbol{\theta}^{\tau,(m)}\|_{\infty}^2 \\ & \lesssim \frac{\mu^2 r^2 \lambda_{\max}^{*2} \log d}{d^2 p} \|\overline{\mathbf{u}}^{\tau,(m)}\|_{\infty}^2. \end{aligned}$$

Apply the Bernstein inequality to yield that with probability at least $1 - O(d^{-11})$,

$$\sum_{j,k \in [d]} T_{m,j,k}^* \overline{u}_j^{\tau,(m)} \theta_k^{\tau,(m)} (p^{-1} \chi_{m,j,k} - 1) \lesssim \sqrt{V_1 \log d} + L_1 \log d \lesssim \frac{\mu r \lambda_{\max}^* \log d}{d \sqrt{p}} \|\overline{\mathbf{u}}^{\tau,(m)}\|_{\infty},$$

where the last inequality holds as long as $p \gg d^{-2} \log^2 d$.

- Regarding α_2 (cf. (223)), it is straightforward to compute that

$$\begin{aligned} L_2 & := \left\| p^{-1} \overline{u}_j^{\tau,(m)} \theta_k^{\tau,(m)} E_{m,j,k} \chi_{m,j,k} \right\|_{\psi_1} \leq \frac{\sigma}{p} \|\overline{\mathbf{u}}^{\tau,(m)}\|_{\infty} \|\boldsymbol{\theta}^{\tau,(m)}\|_{\infty} \\ & \lesssim \frac{\sigma}{p} \sqrt{\frac{\mu r \log d}{d}} \|\overline{\mathbf{u}}^{\tau,(m)}\|_{\infty}, \end{aligned}$$

with $\|\cdot\|_{\psi_1}$ denoting the sub-exponential norm, and

$$\begin{aligned} V_2 & := \mathbb{E} \left[\left(\sum_{j,k \in [d]} \overline{u}_j^{\tau,(m)} \theta_k^{\tau,(m)} p^{-1} E_{m,j,k} \chi_{m,j,k} \right)^2 \right] \\ & = \sum_{j,k \in [d]} (\overline{u}_j^{\tau,(m)})^2 (\theta_k^{\tau,(m)})^2 \mathbb{E}[p^{-2} E_{m,j,k}^2 \chi_{m,j,k}] \\ & \leq \frac{\sigma^2}{p} \|\overline{\mathbf{u}}^{\tau,(m)}\|_2^2 \|\boldsymbol{\theta}^{\tau,(m)}\|_2^2 \lesssim \frac{\sigma^2 r d \log d}{p} \|\overline{\mathbf{u}}^{\tau,(m)}\|_{\infty}^2. \end{aligned}$$

Then the Bernstein inequality reveals that with probability at least $1 - O(d^{-11})$,

$$\left| p^{-1} \sum_{j,k \in [d]} \overline{u}_j^{\tau,(m)} \theta_k^{\tau,(m)} E_{m,j,k} \chi_{m,j,k} \right| \lesssim L_2 \log^2 d + \sqrt{V_2 \log d} \leq \sigma \sqrt{\frac{rd \log^2 d}{p}} \|\overline{\mathbf{u}}^{\tau,(m)}\|_{\infty},$$

where the last inequality follows from our sample size condition.

– Substituting these into (223), we arrive at

$$\left| (\widehat{\mathbf{M}}^{\tau, (m)} - \mathbf{M}^{\tau, (m)})_{m, :} \bar{\mathbf{u}}^{\tau, (m)} \right| \lesssim \left\{ \frac{\mu r \lambda_{\max}^* \log d}{d\sqrt{p}} + \sigma \sqrt{\frac{rd \log^2 d}{p}} \right\} \|\bar{\mathbf{u}}^{\tau, (m)}\|_{\infty}. \quad (224)$$

- For the remaining entries of $(\widehat{\mathbf{M}}^{\tau, (m)} - \mathbf{M}^{\tau, (m)}) \bar{\mathbf{u}}^{\tau, (m)}$, we have

$$\begin{aligned} & (\widehat{\mathbf{M}}^{\tau, (m)} - \mathbf{M}^{\tau, (m)})_{i, :} \bar{\mathbf{u}}^{\tau, (m)} \\ &= \theta_m^{\tau, (m)} \left(\sum_{j: j \neq m} \bar{u}_j^{\tau, (m)} (T_{i, j, m}^* (p^{-1} \chi_{i, j, m} - 1) + p^{-1} E_{i, j, m} \chi_{i, j, m}) \right) \\ & \quad + \bar{v}_m^{\tau, (m)} \left(\sum_{k \in [d]} \theta_k^{\tau, (m)} (T_{i, m, k}^* (p^{-1} \chi_{i, m, k} - 1) + p^{-1} E_{i, m, k} \chi_{i, m, k}) \right) \end{aligned}$$

for any $i \neq m$. From Lemma D.9 and (200), we have, with probability at least $1 - O(d^{-11})$, that

$$\begin{aligned} & (\theta_m^{\tau, (m)})^2 \sum_{i: i \neq m} \left(\sum_{j: j \neq m} T_{i, j, m}^* \bar{u}_j^{\tau, (m)} (p^{-1} \chi_{i, j, m} - 1) \right)^2 + (\bar{v}_m^{\tau, (m)})^2 \sum_{i: i \neq m} \left(\sum_{k \in [d]} T_{i, m, k}^* \theta_k^{\tau, (m)} (p^{-1} \chi_{i, m, k} - 1) \right)^2 \\ & \lesssim \frac{\mu r \lambda_{\max}^{*2}}{dp} \|\boldsymbol{\theta}^{\tau, (m)}\|_{\infty}^2 \|\bar{\mathbf{u}}^{\tau, (m)}\|_{\infty}^2 \\ & \lesssim \frac{\mu^2 r^2 \lambda_{\max}^{*2} \log d}{d^2 p} \|\bar{\mathbf{u}}^{\tau, (m)}\|_{\infty}^2. \end{aligned}$$

Combined with (198) and (200), Lemma D.10 reveals that with probability at least $1 - O(d^{-11})$,

$$\begin{aligned} & (\theta_m^{\tau, (m)})^2 \sum_{i: i \neq m} \left(p^{-1} \sum_{j: j \neq m} \bar{u}_j^{\tau, (m)} E_{i, j, m} \chi_{i, j, m} \right)^2 + (\bar{v}_m^{\tau, (m)})^2 \sum_{i: i \neq m} \left(p^{-1} \sum_{k \in [d]} \theta_k^{\tau, (m)} E_{i, m, k} \chi_{i, m, k} \right)^2 \\ & \lesssim \frac{\sigma^2 d}{p} \|\bar{\mathbf{u}}^{\tau, (m)}\|_2^2 \|\boldsymbol{\theta}^{\tau, (m)}\|_{\infty}^2 + \frac{\sigma^2 d}{p} \|\boldsymbol{\theta}^{\tau, (m)}\|_2^2 \|\bar{\mathbf{u}}^{\tau, (m)}\|_{\infty}^2 + \frac{\sigma^2 \log^5 d}{p^2} \|\bar{\mathbf{u}}^{\tau, (m)}\|_{\infty}^2 \|\boldsymbol{\theta}^{\tau, (m)}\|_{\infty}^2 \\ & \lesssim \frac{\sigma^2 r d \log d}{p} \|\bar{\mathbf{u}}^{\tau, (m)}\|_{\infty}^2, \end{aligned}$$

which implies that

$$\sum_{i: i \neq m} \left((\widehat{\mathbf{M}}^{\tau, (m)} - \mathbf{M}^{\tau, (m)})_{i, :} \bar{\mathbf{u}}^{\tau, (m)} \right)^2 \lesssim \left\{ \frac{\mu^2 r^2 \lambda_{\max}^{*2} \log d}{d^2 p} + \frac{\sigma^2 r d \log d}{p} \right\} \|\bar{\mathbf{u}}^{\tau, (m)}\|_{\infty}^2. \quad (225)$$

- Therefore, combine (224) and (225) to obtain that

$$\|(\widehat{\mathbf{M}}^{\tau, (m)} - \mathbf{M}^{\tau, (m)}) \bar{\mathbf{u}}^{\tau, (m)}\|_2 \lesssim \left\{ \frac{\mu r \lambda_{\max}^* \log d}{d\sqrt{p}} + \sigma \sqrt{\frac{rd \log^2 d}{p}} \right\} \|\bar{\mathbf{u}}^{\tau, (m)}\|_{\infty}.$$

B.8 Proof of Lemma 5.16

By definition, the only possible difference between $\bar{\mathbf{u}}^{\tau}$ and $\boldsymbol{\nu}^{\tau}$ lies in how their global signs are chosen. To show that $\bar{\mathbf{u}}^{\tau} = \boldsymbol{\nu}^{\tau}$, we first claim for the moment that

$$|\langle p^{-1} \mathbf{T}, (\bar{\mathbf{u}}^{\tau})^{\otimes 3} \rangle - \lambda_1^*| \lesssim \mathcal{E}_{\text{proj}} \cdot \lambda_1^*, \quad (226)$$

where $\mathcal{E}_{\text{proj}}$ is defined in (79). Given that $\mathcal{E}_{\text{proj}} \ll 1$ under our sample size, noise and rank condition, this immediately implies that $\langle p^{-1} \mathbf{T}, (\bar{\mathbf{u}}^{\tau})^{\otimes 3} \rangle > 0$. Consequently, by construction, the global signs of $\bar{\mathbf{u}}^{\tau}$ and $\boldsymbol{\nu}^{\tau}$ coincide. Moreover, from (84) and the condition that $\mathcal{E}_{\text{loo}} \sqrt{\mu r \log d/d} \ll 1$, one also knows that $\langle p^{-1} \mathbf{T}^{(m)}, (\bar{\mathbf{u}}^{\tau, (m)})^{\otimes 3} \rangle > 0$ and hence the global signs of $\bar{\mathbf{u}}^{\tau, (m)}$ and $\boldsymbol{\nu}^{\tau, (m)}$ also coincide.

In addition, recall that $\lambda_\tau = \langle p^{-1}\mathbf{T}, (\boldsymbol{\nu}^\tau)^{\otimes 3} \rangle$. One thus has

$$\lambda_\tau = \langle p^{-1}\mathbf{T}, (\boldsymbol{\nu}^\tau)^{\otimes 3} \rangle = \langle p^{-1}\mathbf{T}, (\bar{\mathbf{u}}^\tau)^{\otimes 3} \rangle, \quad (227)$$

which taken collectively with (226) justifies (85).

The rest of the proof then comes down to establishing the claim (226). Towards this, we first decompose

$$\begin{aligned} \langle p^{-1}\mathbf{T}, (\bar{\mathbf{u}}^\tau)^{\otimes 3} \rangle - \lambda_1^* &= \underbrace{\langle p^{-1}\mathbf{T}, (\bar{\mathbf{u}}^\tau)^{\otimes 3} \rangle - \langle \mathbf{T}^*, (\bar{\mathbf{u}}^\tau)^{\otimes 3} \rangle}_{=: \beta_1} \\ &+ \underbrace{\langle \mathbf{T}^*, (\bar{\mathbf{u}}^\tau)^{\otimes 3} \rangle - \langle \mathbf{T}^*, (\bar{\mathbf{u}}_1^*)^{\otimes 3} \rangle}_{=: \beta_2} + \underbrace{\langle \mathbf{T}^*, (\bar{\mathbf{u}}_1^*)^{\otimes 3} \rangle - \lambda_1^*}_{=: \beta_3}. \end{aligned} \quad (228)$$

In what follows, we shall upper bound these three terms separately.

B.8.1 Controlling β_1

Let us start with β_1 . For simplicity of notation, let us define $\boldsymbol{\Delta}_1 := \bar{\mathbf{u}}^\tau - \bar{\mathbf{u}}_1^*$. By construction, \mathbf{T} and \mathbf{T}^* are symmetric. We then can expand

$$\begin{aligned} \langle p^{-1}\mathbf{T} - \mathbf{T}^*, (\bar{\mathbf{u}}^\tau)^{\otimes 3} \rangle &= \langle p^{-1}\mathbf{T} - \mathbf{T}^*, (\bar{\mathbf{u}}_1^* + \boldsymbol{\Delta}_1)^{\otimes 3} \rangle \\ &= \langle p^{-1}\mathbf{T} - \mathbf{T}^*, (\bar{\mathbf{u}}_1^*)^{\otimes 3} \rangle + 3 \langle p^{-1}\mathbf{T} - \mathbf{T}^*, \boldsymbol{\Delta}_1 \otimes (\bar{\mathbf{u}}_1^*)^{\otimes 2} \rangle + 3 \langle p^{-1}\mathbf{T} - \mathbf{T}^*, \boldsymbol{\Delta}_1^{\otimes 2} \otimes \bar{\mathbf{u}}_1^* \rangle \\ &\quad + \langle p^{-1}\mathbf{T} - \mathbf{T}^*, \boldsymbol{\Delta}_1^{\otimes 3} \rangle. \end{aligned} \quad (229)$$

We first look at the first term of (229) which only consists of $\bar{\mathbf{u}}_1^*$. As shown in (173), with probability at least $1 - O(d^{-11})$, one has

$$\|(p^{-1}\mathbf{T} - \mathbf{T}^*) \times_3 \bar{\mathbf{u}}_1^*\| \lesssim \frac{\mu\sqrt{r} \lambda_{\max}^* \sqrt{\log d}}{d\sqrt{p}} + \sigma \sqrt{\frac{d \log d}{p}} \leq \mathcal{E}_{\text{proj}} \lambda_{\min}^*.$$

It follows that

$$|\langle p^{-1}\mathbf{T} - \mathbf{T}^*, (\bar{\mathbf{u}}_1^*)^{\otimes 3} \rangle| \leq \|(p^{-1}\mathbf{T} - \mathbf{T}^*) \times_3 \bar{\mathbf{u}}_1^*\| \|\bar{\mathbf{u}}_1^*\|_2^2 \lesssim \mathcal{E}_{\text{proj}} \lambda_{\min}^*,$$

where we recall the definition of $\mathcal{E}_{\text{proj}}$ in (79). As for the term linear in $\boldsymbol{\Delta}_1$, by Lemma 5.13, we know that $\|\boldsymbol{\Delta}_1\|_2 \lesssim \mathcal{E}_{\text{proj}}$. As a result, one has

$$|\langle p^{-1}\mathbf{T} - \mathbf{T}^*, \boldsymbol{\Delta}_1 \otimes (\bar{\mathbf{u}}_1^*)^{\otimes 2} \rangle| \leq \|(p^{-1}\mathbf{T} - \mathbf{T}^*) \times_3 \bar{\mathbf{u}}_1^*\| \|\bar{\mathbf{u}}_1^*\|_2 \|\boldsymbol{\Delta}_1\|_2 \lesssim \mathcal{E}_{\text{proj}}^2 \lambda_{\min}^*.$$

We then turn to the quadratic terms in $\boldsymbol{\Delta}_1$. Similar to the above arguments, one can deduce that

$$|\langle p^{-1}\mathbf{T} - \mathbf{T}^*, \boldsymbol{\Delta}_1^{\otimes 2} \otimes \bar{\mathbf{u}}_1^* \rangle| \leq \|(p^{-1}\mathbf{T} - \mathbf{T}^*) \times_3 \bar{\mathbf{u}}_1^*\| \|\boldsymbol{\Delta}_1\|_2^2 \lesssim \mathcal{E}_{\text{proj}}^3 \lambda_{\min}^*.$$

Finally, we can simply upper bound the last term in (173) by

$$|\langle p^{-1}\mathbf{T} - \mathbf{T}^*, \boldsymbol{\Delta}_1^{\otimes 3} \rangle| \leq \|p^{-1}\mathbf{T} - \mathbf{T}^*\| \|\boldsymbol{\Delta}_1\|_2^3 \lesssim \mathcal{E}_{\text{op}} \mathcal{E}_{\text{proj}}^3 \lambda_{\min}^* \ll \mathcal{E}_{\text{proj}}^3 \lambda_{\min}^*,$$

where the last step is due to the fact that $\mathcal{E}_{\text{op}} \ll 1$. By our sample size, noise and rank conditions, one has $\mathcal{E}_{\text{proj}} \ll 1$. Putting these bounds together reveals that

$$|\langle p^{-1}\mathbf{T} - \mathbf{T}^*, (\bar{\mathbf{u}}^\tau)^{\otimes 3} \rangle| \lesssim \mathcal{E}_{\text{proj}} \lambda_{\min}^*. \quad (230)$$

B.8.2 Controlling β_2

Recall the definition of β_2 in (228) and $\mathbf{\Delta}_1 = \bar{\mathbf{u}}^\tau - \bar{\mathbf{u}}_1^*$. We can further decompose

$$\langle \mathbf{T}^*, (\bar{\mathbf{u}}^\tau)^{\otimes 3} - (\bar{\mathbf{u}}_1^*)^{\otimes 3} \rangle = 3 \langle \mathbf{T}^*, (\bar{\mathbf{u}}_1^*)^{\otimes 2} \otimes \mathbf{\Delta}_1 \rangle + 3 \langle \mathbf{T}^*, \mathbf{\Delta}_1^{\otimes 2} \otimes \bar{\mathbf{u}}_1^* \rangle + \langle \mathbf{T}^*, \mathbf{\Delta}_1^{\otimes 3} \rangle.$$

We first consider the first term which is linear in $\mathbf{\Delta}_1$. Since \mathbf{T}^* is a symmetric tensor, we have

$$\mathbf{T}^* \times_1 \bar{\mathbf{u}}_1^* = \mathbf{T}^* \times_2 \bar{\mathbf{u}}_1^* = \mathbf{T}^* \times_3 \bar{\mathbf{u}}_1^* = \sum_{s \in [r]} \lambda_1^* \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle \bar{\mathbf{u}}_s^* \bar{\mathbf{u}}_s^{*\top}.$$

By Lemma D.1, one has

$$\|\mathbf{T}^* \times_1 \bar{\mathbf{u}}_1^*\| \leq \max_{1 \leq s \leq r} |\lambda_1^* \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle| \|\bar{\mathbf{U}}^*\|^2 \leq (\lambda_1^* + \lambda_{\max}^* \sqrt{\mu/d}) \lesssim \lambda_1^*,$$

which arises from $\max_{s \neq i} |\langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle| \leq \sqrt{\mu/d}$ and $\|\bar{\mathbf{U}}^*\| \lesssim 1$ as long as $r\sqrt{\mu/d} \ll 1$. As a result, one has

$$|\langle \mathbf{T}^*, (\bar{\mathbf{u}}_1^*)^{\otimes 2} \otimes \mathbf{\Delta}_1 \rangle| \leq \|\mathbf{T}^* \times_1 \bar{\mathbf{u}}_1^*\| \|\bar{\mathbf{u}}_1^*\|_2 \|\mathbf{\Delta}_1\|_2 \lesssim \lambda_1^* \|\mathbf{\Delta}_1\|_2 \lesssim \mathcal{E}_{\text{proj}} \lambda_1^*.$$

In a similar manner, we also know that

$$|\langle \mathbf{T}^*, \mathbf{\Delta}_1^{\otimes 2} \otimes \bar{\mathbf{u}}_1^* \rangle| \leq \|\mathbf{T}^* \times_1 \bar{\mathbf{u}}_1^*\| \|\mathbf{\Delta}_1\|_2^2 \lesssim \mathcal{E}_{\text{proj}}^2 \lambda_1^*.$$

Finally, using the fact that the tensor spectral norm is always less than or equal to that of its matricization, we find that

$$|\langle \mathbf{T}^*, \mathbf{\Delta}_1^{\otimes 3} \rangle| \leq \|\mathbf{T}^*\| \|\mathbf{\Delta}_1\|_2^3 \leq \|\mathbf{A}^*\| \|\mathbf{\Delta}_1\|_2^3 \lesssim \mathcal{E}_{\text{proj}}^3 \lambda_1^*.$$

Combining this with the fact that $\mathcal{E}_{\text{proj}} \ll 1$, we conclude that

$$|\langle \mathbf{T}^*, (\bar{\mathbf{u}}^\tau)^{\otimes 3} - (\bar{\mathbf{u}}_1^*)^{\otimes 3} \rangle| \lesssim \mathcal{E}_{\text{proj}} \lambda_1^*. \quad (231)$$

B.8.3 Controlling β_3

It remains to control β_3 . Straightforward calculation reveals that

$$\langle \mathbf{T}^*, (\bar{\mathbf{u}}_1^*)^{\otimes 3} \rangle = \sum_{s \in [r]} \lambda_s^* \langle (\bar{\mathbf{u}}_s^*)^{\otimes 3}, (\bar{\mathbf{u}}_1^*)^{\otimes 3} \rangle = \lambda_1^* \|\bar{\mathbf{u}}_1^*\|_2^6 + \sum_{s: s \neq 1} \lambda_s^* \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle^6.$$

By the incoherence conditions, we can upper bound

$$|\langle \mathbf{T}^*, (\bar{\mathbf{u}}_1^*)^{\otimes 3} \rangle - \lambda_1^*| = \sum_{s: s \neq 1} \lambda_s^* \langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle^6 \leq r \max_{s: s \neq 1} |\langle \bar{\mathbf{u}}_s^*, \bar{\mathbf{u}}_1^* \rangle|^6 \lambda_{\max}^* \leq \frac{\mu^3 r \lambda_{\max}^*}{d^3}. \quad (232)$$

B.8.4 Combining β_1 , β_2 and β_3

Putting (230), (231) and (232) together, we find that

$$|\langle p^{-1} \mathbf{T}, (\bar{\mathbf{u}}^\tau)^{\otimes 3} \rangle - \lambda_1^*| \lesssim \mathcal{E}_{\text{proj}} \lambda_1^* + \frac{\mu r \lambda_{\max}^*}{d^3} \leq \left(\mathcal{E}_{\text{proj}} + \frac{\mu^3 r}{d^3} \right) \lambda_1^* \asymp \mathcal{E}_{\text{proj}} \lambda_1^*,$$

where the last step follows from the condition $\mu \leq d$, $r \ll \sqrt{d/\mu}$ and the definition $\mathcal{E}_{\text{proj}} \geq \sqrt{\mu r/d}$ (cf. (79)).

B.9 Proof of Lemma 5.17

We first show that for each $i \in [r]$, $(\mathbf{w}_i, \lambda_i)$ and $(\mathbf{w}_i^{(m)}, \lambda_i^{(m)})$ (returned by PRUNE() in Algorithm 3 and Algorithm 6 respectively) satisfy (80), (82) and (84); in other words, we want to show that they correspond to the same index $\tau \in [L]$ (and are hence produced using the same Gaussian random vector \mathbf{g}^τ). The proof idea is this: given that the proposed algorithms select the pair with the largest spectral gap in each round of PRUNE(), it suffices to ensure that there is sufficient separation between the largest and the second largest spectral gaps (so that both algorithms can identify the same τ).

By Lemma 5.12 and union bounds, we know that with probability at least $1 - \delta$, for each $i \in [r]$,

$$\Delta_i^{(1)} - \Delta_i^{(2)} \gtrsim \lambda_{\min}^* / (r\sqrt{\log d}), \quad (233)$$

where we recall that $\Delta_i^\tau := \gamma_i^{\star\tau} - \max_{j:j \neq i} |\gamma_j^{\star\tau}|$, and $\Delta_i^{(1)} \geq \Delta_i^{(2)} \geq \dots \geq \Delta_i^{(L)}$ denote the order statistics of $\{\Delta_i^\tau\}_{\tau=1}^L$ in descending order. As shown in the proof of Lemma 5.13, the spectral gap of \mathbf{M}^τ is well approximated by $\max_i \Delta_i^\tau$, namely,

$$\left| \max_{1 \leq i \leq r} \Delta_i^\tau - (\sigma_1(\mathbf{M}^\tau) - \sigma_2(\mathbf{M}^\tau)) \right| \lesssim \|\mathbf{M}^\tau - \mathbf{M}^{\tau*}\| \ll \frac{\lambda_{\min}^*}{r\sqrt{\log d}}$$

under our sample size, noise and rank conditions (67). Moreover, from Lemma B.2, we see that \mathbf{M}^τ and $\mathbf{M}^{\tau,(m)}$ are extremely close in terms of the spectral norm, i.e.

$$\|\mathbf{M}^\tau - \mathbf{M}^{\tau,(m)}\| \ll \frac{\lambda_{\min}^*}{r\sqrt{\log d}}.$$

This implies that the perturbation incurred by the leave-out-one procedure is relatively small compared to the difference between the largest and the second largest spectral gaps of \mathbf{M}^τ . Consequently, the leave-one-out estimates $\{(\mathbf{w}_i^{(m)}, \lambda_i^{(m)})\}_{i=1}^r$ returned by Algorithm 6 and the true estimates $\{(\mathbf{w}_i, \lambda_i)\}_{i=1}^r$ should correspond to the same trials and should be generated by the same set of Gaussian random vectors. As a result, they obey (69a), (69b) and (69c) for all $1 \leq m \leq d$.

From the discussion above, we also know that as long as $\sigma_1(\mathbf{M}^\tau) - \sigma_2(\mathbf{M}^\tau) \gtrsim \lambda_{\min}^*$, one has $\|\boldsymbol{\nu}^\tau - \bar{\mathbf{u}}_i^*\|_2 \lesssim \mathcal{E}_{\text{proj}}$ for some $i \in [r]$. This is an immediate consequence of Lemma 5.13 and the fact that the spectral gap of \mathbf{M}^τ and $\max_i \Delta_i^\tau$ are extremely close.

It remains to show that our pruning procedure can return estimates of tensor factors without duplicates. Suppose that there exist $1 \leq \tau_1 \neq \tau_2 \leq L$ such that $\|\boldsymbol{\nu}^{\tau_1} - \bar{\mathbf{u}}_i^*\|_2 \lesssim \mathcal{E}_{\text{proj}}$ and $\|\boldsymbol{\nu}^{\tau_2} - \bar{\mathbf{u}}_i^*\|_2 \lesssim \mathcal{E}_{\text{proj}}$ for some $i \in [r]$. By the triangle inequality, one has

$$\begin{aligned} |\langle \boldsymbol{\nu}^{\tau_1}, \boldsymbol{\nu}^{\tau_2} \rangle| &= \left| \|\bar{\mathbf{u}}_i^*\|_2^2 + \langle \boldsymbol{\nu}^{\tau_1} - \bar{\mathbf{u}}_i^*, \bar{\mathbf{u}}_i^* \rangle + \langle \bar{\mathbf{u}}_i^*, \boldsymbol{\nu}^{\tau_2} - \bar{\mathbf{u}}_i^* \rangle + \langle \boldsymbol{\nu}^{\tau_1} - \bar{\mathbf{u}}_i^*, \boldsymbol{\nu}^{\tau_2} - \bar{\mathbf{u}}_i^* \rangle \right| \\ &\geq 1 - \|\boldsymbol{\nu}^{\tau_1} - \bar{\mathbf{u}}_i^*\|_2 - \|\boldsymbol{\nu}^{\tau_2} - \bar{\mathbf{u}}_i^*\|_2 - \|\boldsymbol{\nu}^{\tau_1} - \bar{\mathbf{u}}_i^*\|_2 \|\boldsymbol{\nu}^{\tau_2} - \bar{\mathbf{u}}_i^*\|_2 \\ &\geq 1 - 2\mathcal{E}_{\text{proj}} - \mathcal{E}_{\text{proj}}^2 \geq 1 - 3\mathcal{E}_{\text{proj}}, \end{aligned}$$

provided that $\mathcal{E}_{\text{proj}} \ll 1$. In addition, for any $j \neq i, j \in [r]$, we know that exists some $1 \leq \tau_3 \leq L$ such that

$$\|\boldsymbol{\nu}^{\tau_3} - \bar{\mathbf{u}}_j^*\|_2 \lesssim \mathcal{E}_{\text{proj}}.$$

Recall our incoherence condition in (8c). It is easy to see that

$$\begin{aligned} |\langle \boldsymbol{\nu}^{\tau_1}, \boldsymbol{\nu}^{\tau_3} \rangle| &= |\langle \bar{\mathbf{u}}_i^*, \bar{\mathbf{u}}_j^* \rangle + \langle \boldsymbol{\nu}^{\tau_1} - \bar{\mathbf{u}}_i^*, \bar{\mathbf{u}}_j^* \rangle + \langle \bar{\mathbf{u}}_i^*, \boldsymbol{\nu}^{\tau_3} - \bar{\mathbf{u}}_j^* \rangle + \langle \boldsymbol{\nu}^{\tau_1} - \bar{\mathbf{u}}_i^*, \boldsymbol{\nu}^{\tau_3} - \bar{\mathbf{u}}_j^* \rangle| \\ &\leq |\langle \bar{\mathbf{u}}_i^*, \bar{\mathbf{u}}_j^* \rangle| + \|\boldsymbol{\nu}^{\tau_1} - \bar{\mathbf{u}}_i^*\|_2 + \|\boldsymbol{\nu}^{\tau_3} - \bar{\mathbf{u}}_j^*\|_2 + \|\boldsymbol{\nu}^{\tau_1} - \bar{\mathbf{u}}_i^*\|_2 \|\boldsymbol{\nu}^{\tau_3} - \bar{\mathbf{u}}_j^*\|_2 \\ &\leq \sqrt{\mu/d} + 2\mathcal{E}_{\text{proj}} + \mathcal{E}_{\text{proj}}^2 \leq \sqrt{\mu/d} + 3\mathcal{E}_{\text{proj}} \ll 1 - 3\mathcal{E}_{\text{proj}}, \end{aligned}$$

with the proviso that $\mu \ll d$ and $\mathcal{E}_{\text{proj}} \ll 1$.

The above argument reveals a clear separation between $|\langle \boldsymbol{\nu}^{\tau_1}, \boldsymbol{\nu}^{\tau_3} \rangle|$ and $|\langle \boldsymbol{\nu}^{\tau_1}, \boldsymbol{\nu}^{\tau_2} \rangle|$. As an immediate consequence, the proposed pruning procedure successfully removes all duplication while securing an estimate for each tensor factor.

B.10 Proof of Corollary 5.11

Fix any arbitrary small constant $\delta > 0$. From Theorems 5.9-5.10 and the assumptions of Theorem 2.8, one knows that with probability exceeding $1 - \delta$, there exists a permutation $\pi(\cdot) : [d] \mapsto [d]$ such that for all $1 \leq i \leq r$,

$$\begin{aligned} \|\mathbf{w}^i - \bar{\mathbf{u}}_{\pi(i)}^*\|_2 &\leq \delta, & \|\mathbf{w}^i - \bar{\mathbf{u}}_{\pi(i)}^*\|_\infty &\leq \delta\sqrt{\frac{1}{d}}, & |\lambda_i - \lambda_{\pi(i)}^*| &\leq \delta\lambda_{\max}^*, \\ \|\mathbf{w}^i - \mathbf{w}^{i,(m)}\|_2 &\leq \delta\sqrt{\frac{1}{d}}, & |\lambda_i - \lambda_i^{(m)}| &\leq \delta\sqrt{\frac{1}{d}}\lambda_{\max}^*, & \|(\mathbf{w}^i - \bar{\mathbf{u}}_{\pi(i)}^*)_m\|_2 &\leq \delta\sqrt{\frac{1}{d}} \end{aligned}$$

for some $0 < \delta \ll 1/(\mu^{3/2}r) < 1$. To prove the corollary, we shall just combine the above results.

Without loss of generality, assume that $\pi(i) = i$ for each $i \in [r]$. Given that $\delta \ll 1$ and $\kappa \asymp 1$, by the triangle inequality, one has $\lambda_i \asymp \lambda_i^*$ for all $i \in [r]$, which further implies that

$$|\lambda_i^{1/3} - \lambda_i^{*1/3}| \lesssim \frac{|\lambda_i - \lambda_i^*|}{\lambda_i^{*2/3}} \lesssim \delta\lambda_{\max}^{*1/3}.$$

Consequently, we can apply the triangle inequality to demonstrate that: for each $1 \leq i \leq r$,

$$\|\lambda_i^{1/3}\mathbf{w}^i - \mathbf{u}_i^*\|_2 \leq |\lambda_i^{1/3} - \lambda_i^{*1/3}| \|\bar{\mathbf{u}}_i^*\|_2 + \lambda_i^{*1/3} \|\mathbf{w}^i - \bar{\mathbf{u}}_i^*\|_2 \lesssim \delta\lambda_{\max}^{*1/3}.$$

Arguing similarly, we also see that

$$\begin{aligned} \|\lambda_i^{1/3}\mathbf{w}^i - \mathbf{u}_i^*\|_\infty &\lesssim \delta\sqrt{\frac{1}{d}}\lambda_{\max}^{*1/3}, \\ \|\lambda_i^{1/3}\mathbf{w}^i - (\lambda_i^{(m)})^{1/3}\mathbf{w}^{i,(m)}\|_2 &\lesssim \delta\sqrt{\frac{1}{d}}\lambda_{\max}^{*1/3}, \\ \left| \left((\lambda_i^{(m)})^{1/3}\mathbf{w}^{i,(m)} - \mathbf{u}_i^* \right)_m \right| &\lesssim \delta\sqrt{\frac{1}{d}}\lambda_{\max}^{*1/3} \end{aligned}$$

hold for all $i \in [r]$ and $m \in [d]$. Recall that $\mathbf{U}^0 = [\lambda_i^{1/3}\mathbf{w}^i]_{1 \leq i \leq r}$. One can deduce that

$$\begin{aligned} \|\mathbf{U}^0 - \mathbf{U}^*\|_{\text{F}} &\lesssim \delta\sqrt{r}\lambda_{\max}^{*1/3} \lesssim \delta\|\mathbf{U}^*\|_{\text{F}}, \\ \|\mathbf{U}^0 - \mathbf{U}^*\|_{2,\infty} &\lesssim \delta\sqrt{\frac{r}{d}}\lambda_{\max}^{*1/3} \lesssim \delta\|\mathbf{U}^*\|_{2,\infty}, \\ \|\mathbf{U}^0 - \mathbf{U}^{0,(m)}\|_{2,\infty} &\lesssim \delta\sqrt{\frac{r}{d}}\lambda_{\max}^{*1/3} \lesssim \delta\|\mathbf{U}^*\|_{2,\infty}, \\ \|(\mathbf{U}^{0,(m)} - \mathbf{U}^*)_{m,:}\|_2 &\lesssim \delta\sqrt{\frac{r}{d}}\lambda_{\max}^{*1/3} \lesssim \delta\|\mathbf{U}^*\|_{2,\infty}, \end{aligned}$$

where we have used the condition that $\kappa \asymp 1$ and the fact that $\|\mathbf{U}^*\|_{\text{F}} \geq \sqrt{r}\lambda_{\min}^{*1/3}$ and $\|\mathbf{U}^*\|_{2,\infty} \geq \|\mathbf{U}^*\|_{\text{F}}/\sqrt{d}$.

C Proof of Corollary 2.9

This section establishes Corollary 2.9. First of all, it is easy to see that: given the estimation accuracy established in Theorem 2.8, the permutation matrices that best match \mathbf{U}^t to \mathbf{U}^* remain unchanged as t increases. Therefore, we assume without loss of generality that $\mathbf{I}_r = \arg \min_{\mathbf{\Pi} \in \text{perm}_r} \|\mathbf{U}^t \mathbf{\Pi} - \mathbf{U}^*\|$ for all $t \geq 0$.

Suppose that $r\sqrt{\mu/d} \ll 1$. We claim for the moment that: if a matrix $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_r] \in \mathbb{R}^{d \times r}$ satisfies

$$\|\mathbf{U} - \mathbf{U}^*\|_{\text{F}} \leq \delta\|\mathbf{U}^*\|_{\text{F}} \quad \text{and} \quad \|\mathbf{U} - \mathbf{U}^*\|_{2,\infty} \leq \delta\|\mathbf{U}^*\|_{2,\infty}$$

for any $0 \leq \delta \ll 1/(\mu^{3/2}r) \leq 1$, then one has

$$\|\mathbf{T} - \mathbf{T}^*\|_{\text{F}} \lesssim \delta \|\mathbf{T}^*\|_{\text{F}} \quad \text{and} \quad \|\mathbf{T} - \mathbf{T}^*\|_{\infty} \lesssim \delta \sqrt{\mu^3 r} \|\mathbf{T}^*\|_{\infty}, \quad (234)$$

where $\mathbf{T} := \sum_{i=1}^r \mathbf{u}_i \otimes \mathbf{u}_i \otimes \mathbf{u}_i$. As already shown in the analysis of Theorem 2.8, one has

$$\begin{aligned} \|\mathbf{U}^t - \mathbf{U}^*\|_{\text{F}} &\lesssim \left(C_1 \frac{\rho^{t+1}}{\mu^{3/2}r} + C_2 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{\text{F}}, \\ \|\mathbf{U}^t - \mathbf{U}^*\|_{2,\infty} &\lesssim \left(C_3 \rho^{t+1} \frac{\rho^{t+1}}{\mu^{3/2}r} + C_4 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \|\mathbf{U}^*\|_{2,\infty}, \end{aligned}$$

from which Corollary 2.9 follows immediately.

It remains to prove the claim (234). For notational convenience, let us define $\mathbf{\Delta} := \mathbf{U} - \mathbf{U}^*$ and $\mathbf{\Delta}_s := \mathbf{u}_s - \mathbf{u}_s^*$ for each $1 \leq s \leq r$. Then we can expand

$$\begin{aligned} \mathbf{T} - \mathbf{T}^* &= \sum_{1 \leq s \leq r} \mathbf{u}_s^{\otimes 3} - \mathbf{u}_s^{*\otimes 3} = \sum_{1 \leq s \leq r} \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{\Delta}_s + \sum_{1 \leq s \leq r} \mathbf{u}_s^* \otimes \mathbf{\Delta}_s \otimes \mathbf{u}_s^* + \sum_{1 \leq s \leq r} \mathbf{\Delta}_s \otimes \mathbf{u}_s^{*\otimes 2} \\ &\quad + \sum_{1 \leq s \leq r} \mathbf{u}_s^* \otimes \mathbf{\Delta}_s^{\otimes 2} + \sum_{1 \leq s \leq r} \mathbf{\Delta}_s \otimes \mathbf{u}_s^* \otimes \mathbf{\Delta}_s + \sum_{1 \leq s \leq r} \mathbf{\Delta}_s^{\otimes 2} \otimes \mathbf{u}_s^* + \sum_{1 \leq s \leq r} \mathbf{\Delta}_s^{\otimes 3}. \end{aligned} \quad (235)$$

(1) Euclidean loss. We first look at the loss measured by $\|\cdot\|_{\text{F}}$. In view of the symmetric structure of tensors, it suffices to control $\sum_{1 \leq s \leq r} \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{\Delta}_s$, $\sum_{1 \leq s \leq r} \mathbf{u}_s^* \otimes \mathbf{\Delta}_s^{\otimes 2}$ and $\sum_{1 \leq s \leq r} \mathbf{\Delta}_s^{\otimes 3}$.

Let us define $\mathbf{W}_1 := [\mathbf{u}_s^* \otimes \mathbf{\Delta}_s]_{1 \leq s \leq r} \in \mathbb{R}^{d^2 \times r}$ and $\mathbf{W}_2 := [\mathbf{\Delta}_s^{\otimes 2}]_{1 \leq s \leq r} \in \mathbb{R}^{d^2 \times r}$. Recalling the fact that $\|\mathbf{U}^*\| \leq \|\overline{\mathbf{U}}^*\| \lambda_{\max}^{*1/3} \lesssim \lambda_{\max}^{*1/3}$ (established in Lemma D.1), we have

$$\begin{aligned} \left\| \sum_{1 \leq s \leq r} \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{\Delta}_s \right\|_{\text{F}} &= \|\mathbf{U}^* \mathbf{W}_1^{\top}\|_{\text{F}} \leq \|\mathbf{U}^*\| \|\mathbf{W}_1\|_{\text{F}} \lesssim \lambda_{\max}^{*1/3} \|\mathbf{W}_1\|_{\text{F}}, \\ \left\| \sum_{1 \leq s \leq r} \mathbf{u}_s^* \otimes \mathbf{\Delta}_s^{\otimes 2} \right\|_{\text{F}} &= \|\mathbf{U}^* \mathbf{W}_2^{\top}\|_{\text{F}} \leq \|\mathbf{U}^*\| \|\mathbf{W}_2\|_{\text{F}} \lesssim \lambda_{\max}^{*1/3} \|\mathbf{W}_2\|_{\text{F}}, \\ \left\| \sum_{1 \leq s \leq r} \mathbf{\Delta}_s^{\otimes 3} \right\|_{\text{F}} &= \|\mathbf{\Delta} \mathbf{W}_2^{\top}\|_{\text{F}} \leq \|\mathbf{\Delta}\| \|\mathbf{W}_2\|_{\text{F}} \leq \|\mathbf{\Delta}\|_{\text{F}} \|\mathbf{W}_2\|_{\text{F}}. \end{aligned}$$

It then suffices to control $\|\mathbf{W}_1\|_{\text{F}}$ and $\|\mathbf{W}_2\|_{\text{F}}$. If $\|\mathbf{\Delta}\|_{\text{F}} \leq \delta \|\mathbf{U}^*\|_{\text{F}} \leq \delta \sqrt{r} \lambda_{\max}^{*1/3}$, then it is easy to bound

$$\begin{aligned} \|\mathbf{W}_1\|_{\text{F}}^2 &= \sum_{1 \leq s \leq r} \|\mathbf{u}_s^* \otimes \mathbf{\Delta}_s\|_2^2 = \sum_{1 \leq s \leq r} \|\mathbf{u}_s^*\|_2^2 \|\mathbf{\Delta}_s\|_2^2 \leq \max_{1 \leq s \leq r} \|\mathbf{u}_s^*\|_2^2 \|\mathbf{\Delta}\|_{\text{F}}^2 \leq \delta^2 r \lambda_{\max}^{*4/3}, \\ \|\mathbf{W}_2\|_{\text{F}}^2 &= \sum_{1 \leq s \leq r} \|\mathbf{\Delta}_s^{\otimes 2}\|_2^2 = \sum_{1 \leq s \leq r} \|\mathbf{\Delta}_s\|_2^4 \leq \max_{1 \leq s \leq r} \|\mathbf{\Delta}_s\|_2^2 \|\mathbf{\Delta}\|_{\text{F}}^2 \leq \|\mathbf{\Delta}\|_{\text{F}}^4 \leq \delta^4 r^2 \lambda_{\max}^{*4/3}. \end{aligned}$$

Therefore, one has

$$\left\| \sum_{1 \leq s \leq r} \mathbf{u}_s^{*\otimes 2} \otimes \mathbf{\Delta}_s \right\|_{\text{F}} \lesssim \delta \sqrt{r} \lambda_{\max}^*, \quad (236)$$

$$\left\| \sum_{1 \leq s \leq r} \mathbf{u}_s^* \otimes \mathbf{\Delta}_s^{\otimes 2} \right\|_{\text{F}} \lesssim \delta^2 r \lambda_{\max}^*, \quad (237)$$

$$\left\| \sum_{1 \leq s \leq r} \mathbf{\Delta}_s^{\otimes 3} \right\|_{\text{F}} \lesssim \delta^3 r^{3/2} \lambda_{\max}^*. \quad (238)$$

Since $0 \leq \delta \ll r^{-1} \leq 1$, combining (236), (237) and (238) with the fact that $\|\mathbf{T}^*\|_{\text{F}} \geq \sqrt{r} \lambda_{\min}^*/2$ (established in Lemma D.1), we conclude that

$$\|\mathbf{T} - \mathbf{T}^*\| \lesssim \delta \sqrt{r} \lambda_{\max}^* \lesssim \delta \|\mathbf{T}^*\|_{\text{F}}. \quad (239)$$

(2) ℓ_∞ loss. Next, we turn to the $\|\cdot\|_\infty$ loss. Again, it suffices to focus on $\sum_{1 \leq s \leq r} \mathbf{u}_s^{*\otimes 2} \otimes \Delta_s$, $\sum_{1 \leq s \leq r} \mathbf{u}_s^* \otimes \Delta_s^{\otimes 2}$ and $\sum_{1 \leq s \leq r} \Delta_s^{\otimes 3}$. From (104), (105) and (106) shown in the proof of Lemma 5.1, one has

$$\left\| \sum_{1 \leq s \leq r} (\mathbf{u}_s^*)^{\otimes 2} \otimes \Delta_s \right\|_\infty \leq \delta \max_{1 \leq s \leq r} \|\mathbf{u}_s^*\|_\infty \|\mathbf{U}^*\|_{2,\infty}^2 \leq \delta \frac{\mu^{3/2} r \lambda_{\max}^*}{d^{3/2}}, \quad (240)$$

$$\left\| \sum_{1 \leq s \leq r} \mathbf{u}_s^* \otimes \Delta_s^{\otimes 2} \right\|_\infty \leq \delta^2 \max_{1 \leq s \leq r} \|\mathbf{u}_s^*\|_\infty \|\mathbf{U}^*\|_{2,\infty}^2 \leq \frac{\delta^2 \mu^{3/2} r \lambda_{\max}^*}{d^{3/2}}, \quad (241)$$

$$\left\| \sum_{1 \leq s \leq r} \Delta_s^{\otimes 3} \right\|_\infty \leq \delta^3 \|\mathbf{U}^*\|_{2,\infty}^3 \leq \frac{\delta^3 \mu^{3/2} r^{3/2} \lambda_{\max}^*}{d^{3/2}}. \quad (242)$$

Putting (240), (241) and (242) together with the condition that $0 < \delta \ll r^{-1} \leq 1$, we arrive at

$$\|\mathbf{T} - \mathbf{T}^*\|_\infty \lesssim \frac{\delta \mu^{3/2} r \lambda_{\max}^*}{d^{3/2}}.$$

In addition, from the lower bound on $\|\mathbf{T}^*\|_{\text{F}}$, one has

$$\|\mathbf{T}^*\|_\infty \geq \frac{1}{d^{3/2}} \|\mathbf{T}^*\|_{\text{F}} \gtrsim \sqrt{\frac{r}{d^3}} \lambda_{\min}^*,$$

which allows us to conclude that

$$\|\mathbf{T} - \mathbf{T}^*\|_\infty \lesssim \sqrt{\mu^3 r} \delta \|\mathbf{T}^*\|_\infty.$$

D Auxiliary lemmas

This section gathers several auxiliary lemmas that prove useful when establishing our main results.

D.1 Statements of auxiliary lemmas

We begin by stating all auxiliary lemmas formally, with the proofs postponed to subsequent subsections. We shall define

$$\bar{\mathbf{U}}^* := [\bar{\mathbf{u}}_1^*, \dots, \bar{\mathbf{u}}_r^*], \quad \text{with} \quad \bar{\mathbf{u}}_i^* := \mathbf{u}_i^* / \|\mathbf{u}_i^*\|_2. \quad (243)$$

Lemma D.1. *Suppose that Assumption 2.1 holds, and assume that $r\sqrt{\mu/d} \leq c_1$ for some sufficiently small universal constant $c_2 > 0$. Then for d sufficiently large, the matrices \mathbf{A}^* , \mathbf{B}^* and $\mathbf{U}_{\text{orth}}^*$ (defined respectively in (24), (63) and (54)) obey*

$$\begin{aligned} \frac{1}{2} \lambda_{\min}^* &\leq \|\mathbf{A}^*\|_{\text{F}} \leq 2\lambda_{\max}^*, \quad \|\mathbf{A}^*\|_\infty \leq \frac{\sqrt{2\mu r} \lambda_{\max}^*}{d^{3/2}}, \quad \|\mathbf{A}^*\|_{2,\infty} \leq \sqrt{\frac{2\mu r}{d}} \lambda_{\max}^*, \quad \|\mathbf{A}^{*\top}\|_{2,\infty} \leq \frac{\mu\sqrt{2r} \lambda_{\max}^*}{d}, \\ \|\mathbf{A}^*\| &= \lambda_{\max}^* \left(1 + O\left(r\sqrt{\frac{\mu}{d}}\right)\right), \quad \lambda_i(\mathbf{B}^*) = \lambda_{(i)}^{*2} \left(1 + O\left(r\sqrt{\frac{\mu}{d}}\right)\right), \quad i \in [r], \\ \|\mathbf{B}^*\|_{2,\infty} &\leq 2\sqrt{\frac{\mu r}{d}} \lambda_{\max}^{*2}, \quad \|\mathbf{U}_{\text{orth}}^*\|_{2,\infty} \leq \sqrt{\frac{2\mu r}{d}}, \quad \|\bar{\mathbf{U}}^{*\top} \bar{\mathbf{U}}^* - \mathbf{I}\| \leq r\sqrt{\frac{\mu}{d}}. \end{aligned}$$

Here, $\|\mathbf{A}\|_{2,\infty} := \max_i \|\mathbf{A}_{i,:}\|_2$, $\lambda_{(i)}^*$ stands for the i -th largest value in $\{\lambda_i^*\}_{1 \leq i \leq r}$ (or equivalently $\{\|\mathbf{u}_i^*\|_2^3\}_{1 \leq i \leq r}$), and $\lambda_i(\mathbf{B}^*)$ represents the i -th largest eigenvalue of \mathbf{B}^* .

Proof. See Appendix D.2. ■

Lemma D.2. *Let $\mathbf{R} \in \mathbb{R}^{d \times d \times d}$ be a random order-3 tensor with independent entries $\{R_{i,j,k}\}_{i,j,k \in [d]}$ obeying*

$$\mathbb{E}[R_{i,j,k}] = 0, \quad \max_{i,j,k \in [d]} |R_{i,j,k}| \leq B.$$

Define

$$\sigma_{\text{mode}}^2 := \max_{j,k \in [d]} \sum_{i \in [d]} \mathbb{E}[R_{i,j,k}^2] + \max_{i,k \in [d]} \sum_{j \in [d]} \mathbb{E}[R_{i,j,k}^2] + \max_{i,j \in [d]} \sum_{k \in [d]} \mathbb{E}[R_{i,j,k}^2]. \quad (244)$$

Then with probability exceeding $1 - O(d^{-10})$, one has

$$\|\mathbf{R}\| \lesssim B \log^3 d + \sigma_{\text{mode}} \log^{5/2} d. \quad (245)$$

Proof. See Appendix D.3. ■

An immediate consequence of this lemma is the following:

Corollary D.3. *With probability at least $1 - O(d^{-10})$, one has*

$$\|p^{-1} \mathcal{P}_{\Omega}(\mathbf{T}^*) - \mathbf{T}^*\| \lesssim \frac{\sqrt{\mu r} \lambda_{\max}^* \log^3 d}{d^{3/2} p} + \frac{\mu \sqrt{r} \lambda_{\max}^* \log^{5/2} d}{d \sqrt{p}}; \quad (246)$$

$$\|\mathcal{P}_{\Omega}(\mathbf{E})\| \lesssim \sigma(\log^{7/2} d + \sqrt{dp} \log^{5/2} d). \quad (247)$$

Proof. See Appendix D.3. ■

Lemma D.4. *Suppose that $p \gtrsim d^{-2} \log^3 d$ and that $\mu \log^2 d \lesssim d$. Then for any fixed vector $\mathbf{w} \in \mathbb{R}^d$, with probability $1 - O(d^{-10})$, one has*

$$\|(p^{-1} \mathbf{T} - \mathbf{T}^*) \times_3 \mathbf{w}\| \lesssim \|\mathbf{w}\|_{\infty} \sqrt{\frac{\mu r \log d}{dp} \lambda_{\max}^*} + \|\mathbf{w}\|_{\infty} \frac{\sigma \log^{5/2} d}{p} + \|\mathbf{w}\|_2 \sigma \sqrt{\frac{d \log d}{p}},$$

where \times_3 is defined in Section 2.4. The results also holds if we replace \times_3 with \times_1 or \times_2 .

Proof. See Appendix D.4. ■

Lemma D.5. *Let $\{X_{i,j}\}_{1 \leq i \leq r, 1 \leq j \leq L}$ be a sequence of i.i.d. standard Gaussian random variables, where $r \geq 2$ and $L \geq 1$. Consider some quantities $\kappa \geq 1, \Delta > 0, 0 < \delta < 1/2$. There exists some universal constant $C > 0$ such that if*

$$L \geq Cr^{2\kappa^2} (\kappa \sqrt{\log r} + \Delta) \exp(\Delta^2) \log \frac{1}{\delta},$$

then with probability at least $1 - \delta$, there exists some $1 \leq j_0 \leq L$ such that

$$X_{1,j_0} > \kappa \max_{i: 1 < i \leq r} |X_{i,j_0}| + \Delta.$$

In addition, define $\Delta_j := X_{1,j} - \kappa \max_{i: 1 < i \leq r} |X_{i,j}|$ for each $1 \leq j \leq L$. Then with probability at least $1 - 2\delta$,

$$\Delta_{(1)} - \Delta_{(2)} \gtrsim \frac{\delta}{\sqrt{\log L} + \sqrt{\log(1/\delta)}}.$$

where $\Delta_{(1)} \geq \Delta_{(2)} \geq \dots \geq \Delta_{(L)}$ denote the order statistics of $\{\Delta_j\}_{j=1}^L$ in descending order.

Proof. See Appendix D.5. ■

Lemma D.6. *Let \mathbf{U} (resp. \mathbf{V}) be a $d \times r$ matrix with orthonormal columns. Suppose that $\|\mathbf{U}\mathbf{U}^{\top} - \mathbf{V}\mathbf{V}^{\top}\| \leq \delta$. Then for any unit vector $\mathbf{u}_0 \in \mathbb{R}^d$ lying in $\text{span}(\mathbf{U})$, we have*

$$\|\mathcal{P}_{\mathbf{V}}(\mathbf{u}_0)\|_2 \geq \sqrt{1 - \delta^2} \quad \text{and} \quad \|\mathcal{P}_{\mathbf{V}^{\perp}}(\mathbf{u}_0)\|_2 \leq \delta, \quad (248)$$

where we denote by $\mathcal{P}_{\mathbf{V}}(\mathbf{u}_0) := \mathbf{V}\mathbf{V}^{\top} \mathbf{u}_0$ and $\mathcal{P}_{\mathbf{V}^{\perp}}(\mathbf{u}_0) = (\mathbf{I}_d - \mathbf{V}\mathbf{V}^{\top}) \mathbf{u}_0$.

Proof. See Appendix D.6. ■

Additionally, we record several facts concerning the set of Bernoulli random variables $\{\chi_{i,j,k}\}_{1 \leq i,j,k \leq d}$. We recall that

$$\chi_{i,j,k} := \mathbb{1}\{(i,j,k) \in \Omega\}, \quad (249)$$

which is a Bernoulli random variable with mean p .

Lemma D.7. *Suppose that $p \gtrsim d^{-2} \log d$. With probability exceeding $1 - O(d^{-10})$, one has*

$$\sum_{i,j \in [d]} T_{i,j,k}^{*2} (p^{-1} \chi_{i,j,k} - 1)^2 \lesssim \frac{\mu r \lambda_{\max}^{*2}}{dp}, \quad 1 \leq k \leq d. \quad (250)$$

Proof. See Appendix D.7. ■

Lemma D.8. *Suppose that $p \gtrsim d^{-2} \log^2 d$. With probability exceeding $1 - O(d^{-10})$, one has*

$$\sum_{i,j \in [d]} (p^{-1} E_{i,j,k} \chi_{i,j,k})^2 \lesssim \sigma^2 d^2 / p, \quad 1 \leq k \leq d. \quad (251)$$

Proof. See Appendix D.8. ■

Lemma D.9. *Suppose $p \gtrsim d^{-2} \log d$ and $\mu \log^2 d \lesssim d$. Consider any fixed vector $\mathbf{w} \in \mathbb{R}^d$ and any $1 \leq i \leq d$. With probability exceeding $1 - O(d^{-10})$, one has*

$$\sum_{j \in [d]} \left(\sum_{k \in [d]} T_{i,j,k}^* w_k (p^{-1} \chi_{i,j,k} - 1) \right)^2 \lesssim \frac{\mu r \lambda_{\max}^{*2}}{dp} \|\mathbf{w}\|_{\infty}^2. \quad (252)$$

Proof. See Appendix D.9. ■

Lemma D.10. *Consider any fixed vector $\mathbf{w} \in \mathbb{R}^d$. With probability $1 - O(d^{-10})$ one has*

$$\sum_{j \in [d]} \left(\sum_{k \in [d]} w_k E_{i,j,k} \chi_{i,j,k} \right)^2 \lesssim \sigma^2 dp \|\mathbf{w}\|_2^2 + \sigma^2 \|\mathbf{w}\|_{\infty}^2 \log^5 d, \quad 1 \leq i \leq d. \quad (253)$$

Proof. See Appendix D.10. ■

D.2 Proof of Lemma D.1

1. To begin with, the incoherence condition (8b) gives

$$\begin{aligned} \|\mathbf{T}^*\|_{\mathbb{F}}^2 &= \left\langle \sum_{i \in [r]} \mathbf{u}_i^{*\otimes 3}, \sum_{i \in [r]} \mathbf{u}_i^{*\otimes 3} \right\rangle \\ &= \sum_{1 \leq i \leq r} \|\mathbf{u}_i^{*\otimes 3}\|_{\mathbb{F}}^2 + \sum_{1 \leq i \neq j \leq r} \langle \mathbf{u}_i^{*\otimes 3}, \mathbf{u}_j^{*\otimes 3} \rangle \\ &\stackrel{(i)}{=} \sum_{1 \leq i \leq r} \|\mathbf{u}_i^*\|_2^6 + \sum_{1 \leq i \neq j \leq r} \langle \mathbf{u}_i^*, \mathbf{u}_j^* \rangle^3 \\ &\leq r \max_{1 \leq i \leq r} \|\mathbf{u}_i^*\|_2^6 + r^2 \max_{1 \leq i \neq j \leq r} |\langle \mathbf{u}_i^*, \mathbf{u}_j^* \rangle^3| \\ &\leq r \lambda_{\max}^{*2} + r^2 \left(\frac{\mu}{d} \right)^{3/2} \lambda_{\max}^{*2} \\ &\stackrel{(ii)}{\leq} 2r \lambda_{\max}^{*2}, \end{aligned}$$

where we use the fact that $\langle \mathbf{u}^{\otimes 3}, \mathbf{v}^{\otimes 3} \rangle = \langle \mathbf{u}, \mathbf{v} \rangle^3$ in (i), and (ii) arises due to the condition that $3r \leq \sqrt{d/\mu} \leq (d/\mu)^{3/2}$. Using a similar argument, we also know that $\|\mathbf{T}^*\|_{\mathbb{F}}^2 \geq r \lambda_{\min}^{*2}/2$. This combined with the incoherence condition in (8a) yields

$$\|\mathbf{A}^*\|_{\infty} = \|\mathbf{T}^*\|_{\infty} \leq \sqrt{\frac{\mu}{d^3}} \|\mathbf{T}^*\|_{\mathbb{F}} \leq \frac{\sqrt{2\mu r}}{d^{3/2}} \lambda_{\max}^*.$$

2. For any $1 \leq i \leq d$, the ℓ_2 norm of the i -th row of \mathbf{A}^* can be bounded by

$$\begin{aligned}
\|\mathbf{A}_{i,:}^*\|_2^2 &= \left\| \sum_{1 \leq s \leq r} (\mathbf{u}_s^*)_i (\mathbf{u}_s^* \otimes \mathbf{u}_s^*)^\top \right\|_2^2 \\
&= \sum_{1 \leq s \leq r} (\mathbf{u}_s^*)_i^2 \|\mathbf{u}_s^* \otimes \mathbf{u}_s^*\|_2^2 + \sum_{1 \leq s_1 \neq s_2 \leq r} (\mathbf{u}_{s_1}^*)_i (\mathbf{u}_{s_2}^*)_i \langle \mathbf{u}_{s_1}^* \otimes \mathbf{u}_{s_1}^*, \mathbf{u}_{s_2}^* \otimes \mathbf{u}_{s_2}^* \rangle \\
&\leq r \max_{s \in [r]} \|\mathbf{u}_s^*\|_\infty^2 \max_{s \in [r]} \|\mathbf{u}_s^*\|_2^4 + r^2 \max_{1 \leq s \leq r} \|\mathbf{u}_s^*\|_\infty^2 \max_{s_1 \neq s_2} \langle \mathbf{u}_{s_1}^*, \mathbf{u}_{s_2}^* \rangle^2 \\
&\leq \lambda_{\max}^{*2} \left(\frac{\mu r}{d} + \frac{\mu^2 r^2}{d^2} \right) \leq \frac{2\mu r \lambda_{\max}^{*2}}{d},
\end{aligned} \tag{254}$$

where the first inequality follows from the fact that $\langle \mathbf{u} \otimes \mathbf{u}, \mathbf{v} \otimes \mathbf{v} \rangle = \langle \mathbf{u}, \mathbf{v} \rangle^2$, the second inequality holds true due to (8b) and (8c), and the last inequality holds as long as $r \leq d/\mu$. This immediately yields the advertised bound on $\|\mathbf{A}^*\|_{2,\infty}$.

3. For any $\underline{j} \in [d]^2$ (which corresponds to (j_1, j_2)), the ℓ_2 norm of the \underline{j} -th column of \mathbf{A}^* can be upper bounded similarly by

$$\begin{aligned}
\|\mathbf{A}_{:, \underline{j}}^*\|_2^2 &= \left\| \sum_{1 \leq s \leq r} (\mathbf{u}_s^* \otimes \mathbf{u}_s^*)_{\underline{j}} \mathbf{u}_s^* \right\|_2^2 \\
&= \sum_{1 \leq s \leq r} (\mathbf{u}_s^*)_{j_1}^2 (\mathbf{u}_s^*)_{j_2}^2 \|\mathbf{u}_s^*\|_2^2 + \sum_{s_1 \neq s_2} (\mathbf{u}_{s_1}^*)_{j_1} (\mathbf{u}_{s_1}^*)_{j_2} (\mathbf{u}_{s_2}^*)_{j_1} (\mathbf{u}_{s_2}^*)_{j_2} \langle \mathbf{u}_{s_1}^*, \mathbf{u}_{s_2}^* \rangle \\
&\leq r \max_{1 \leq s \leq r} \|\mathbf{u}_s^*\|_\infty^4 \max_{s \in [r]} \|\mathbf{u}_s^*\|_2^2 + r^2 \max_{1 \leq s \leq r} \|\mathbf{u}_s^*\|_\infty^4 \max_{1 \leq s_1 \neq s_2 \leq r} |\langle \mathbf{u}_{s_1}^*, \mathbf{u}_{s_2}^* \rangle| \\
&\leq \lambda_{\max}^{*2} \left(\frac{\mu^2 r}{d^2} + \frac{\mu^{5/2} r^2}{d^{5/2}} \right) \leq \frac{2\mu^2 r \lambda_{\max}^{*2}}{d^2},
\end{aligned} \tag{255}$$

where the second inequality is valid due to (8b) and (8c), and the last inequality holds as long as $r \leq \sqrt{d/\mu}$. This yields the claimed bound regarding $\|\mathbf{A}^{*\top}\|_{2,\infty}$.

4. Regarding the spectrum of \mathbf{A}^* , \mathbf{B}^* , $\mathbf{U}_{\text{orth}}^*$ and $\overline{\mathbf{U}}^*$, we refer the reader to the proof of [CLC⁺20, Corollary 1].

5. We now move on to $\|\mathbf{B}^*\|_{2,\infty}$. For any $i \in [d]$, it is seen that

$$\begin{aligned}
\|\mathbf{B}_{i,:}^*\|_2^2 &= \sum_{j \in [d]} (\mathbf{A}_{i,:}^* \mathbf{A}_{j,:}^{*\top})^2 = \mathbf{A}_{i,:}^* \left(\sum_{j \in [d]} \mathbf{A}_{j,:}^{*\top} \mathbf{A}_{j,:}^* \right) \mathbf{A}_{i,:}^{*\top} = \mathbf{A}_{i,:}^* (\mathbf{A}^{*\top} \mathbf{A}^*) \mathbf{A}_{i,:}^{*\top} \\
&\leq \|\mathbf{A}^{*\top} \mathbf{A}^*\| \|\mathbf{A}_{i,:}^*\|_2^2 \leq 2\lambda_{\max}^{*2} \|\mathbf{A}_{i,:}^*\|_2^2.
\end{aligned}$$

Here, the last line makes use of the bound $\|\mathbf{A}^*\| \leq \lambda_{\max}^* (1 + O(r\sqrt{\mu/d})) \leq 2\lambda_{\max}^*$, which holds if $r\sqrt{\mu/d} \leq c_1$ for some sufficiently small constant $c_1 > 0$. It then follows from (254) that

$$\|\mathbf{B}^*\|_{2,\infty} \leq \sqrt{2} \lambda_{\max}^* \|\mathbf{A}^*\|_{2,\infty} \leq 2\lambda_{\max}^{*2} \sqrt{\frac{\mu r}{d}}. \tag{256}$$

D.3 Proof of Lemma D.2 and Corollary D.3

D.3.1 Proof of Lemma D.2

We start by making the following simple observation: the tensor spectral norm is a 1-Lipschitz function (w.r.t. the Frobenius norm) of the entries of the tensor. This follows since $|\|\mathbf{T}\| - \|\mathbf{R}\|| \leq \|\mathbf{T} - \mathbf{R}\| \leq \|\mathbf{T} - \mathbf{R}\|_F$ holds for any tensor $\mathbf{T}, \mathbf{R} \in \mathbb{R}^{d \times d \times d}$. This allows us to invoke standard concentration results regarding functions of independent random variables.

We shall first develop an upper bound on the mean $\mathbb{E}[\|\mathcal{P}_\Omega(\mathbf{R})\|]$. In view of [NDT15, Corollary 4] and Jensen's inequality, one has

$$\begin{aligned} \mathbb{E}[\|\mathbf{R}\|] &\leq \sqrt{\mathbb{E}[\|\mathbf{R}\|^2]} \\ &\lesssim \left(\mathbb{E} \left[\max_{j,k \in [d]} \sum_{i \in [d]} R_{i,j,k}^2 + \max_{i,j \in [d]} \sum_{k \in [d]} R_{i,j,k}^2 + \max_{i,k \in [d]} \sum_{j \in [d]} R_{i,j,k}^2 \right] \right)^{1/2} \log^{5/2} d. \end{aligned} \quad (257)$$

We then need to bound the quantity presented in (257).

For some $\beta > 0$ to be specified later, one can upper bound

$$\begin{aligned} \mathbb{E} \left[\max_{j,k \in [d]} \sum_{i \in [d]} R_{i,j,k}^2 \right] &= \int_0^\infty \mathbb{P} \left\{ \max_{j,k \in [d]} \sum_{i \in [d]} R_{i,j,k}^2 > t \right\} dt \\ &\leq \beta + \int_\beta^\infty \mathbb{P} \left\{ \max_{j,k \in [d]} \sum_{i \in [d]} R_{i,j,k}^2 > t \right\} dt \\ &\leq \beta + d^2 \int_\beta^\infty \mathbb{P} \left\{ \sum_{i \in [d]} R_{i,j,k}^2 > t \right\} dt. \end{aligned} \quad (258)$$

We shall resort to the Bernstein inequality to bound $\mathbb{P} \left\{ \sum_{i \in [d]} R_{i,j,k}^2 > t \right\}$. It is straightforward to compute that

$$\begin{aligned} M &:= \sum_{i \in [d]} \mathbb{E} [R_{i,j,k}^2] \leq \sigma_{\text{mode}}^2, \\ L &:= \max_{i \in [d]} |R_{i,j,k}^2| \leq B^2, \\ S^2 &:= \sum_{i \in [d]} \mathbb{E} [R_{i,j,k}^4] \leq B^2 \sigma_{\text{mode}}^2. \end{aligned}$$

The Bernstein inequality then tells us that

$$\mathbb{P} \left\{ \sum_{i \in [d]} R_{i,j,k}^2 - M > t \right\} \leq \exp \left(-\frac{3}{8} \min \left\{ \frac{t^2}{S^2}, \frac{t}{L} \right\} \right), \quad t > 0. \quad (259)$$

In particular, this implies that with probability exceeding $1 - O(d^{-20})$,

$$\begin{aligned} \sum_{i \in [d]} R_{i,j,k}^2 &\lesssim M + L \log d + S \sqrt{\log d} \lesssim \sigma_{\text{mode}}^2 + B^2 \log d + \sqrt{B^2 \sigma_{\text{mode}}^2 \log d} \\ &\asymp \sigma_{\text{mode}}^2 + B^2 \log d, \end{aligned}$$

where we have used the AM-GM inequality in the last step. Therefore, by taking

$$\beta := C (\sigma_{\text{mode}}^2 + B^2 \log d)^{1/2}$$

for some sufficiently large constant $C > 0$, we arrive at

$$\beta \geq \frac{C}{3} (M + L \log d + S \sqrt{\log d}) \gg M + L \log d + S \sqrt{\log d}. \quad (260)$$

Given that $\beta \gg M$, for any $t \geq \beta$ one has the following relations about several events

$$\left\{ \sum_{i \in [d]} R_{i,j,k}^2 > t \right\} = \left\{ \sum_{i \in [d]} R_{i,j,k}^2 - M > t - M \right\} \subset \left\{ \sum_{i \in [d]} R_{i,j,k}^2 - M > t - \beta/2 \right\} \quad (261)$$

In addition, it is easily seen that

$$\min \{t^2/S^2, t/L\} \geq \frac{t}{\max \{S/\sqrt{\log d}, L\}} \quad (262)$$

for any $t \geq \beta$ (with β obeying (260)). As a result, one can bound

$$\begin{aligned}
\int_{\beta}^{\infty} \mathbb{P}\left\{\sum_{i \in [d]} R_{i,j,k}^2 > t\right\} dt &\stackrel{(i)}{\leq} \int_{\beta}^{\infty} \mathbb{P}\left\{\sum_{i \in [d]} R_{i,j,k}^2 - M > t - \beta/2\right\} dt \\
&= \int_{\beta/2}^{\infty} \mathbb{P}\left\{\sum_{i \in [d]} R_{i,j,k}^2 - M > t\right\} dt \\
&\stackrel{(ii)}{\leq} \int_{\beta/2}^{\infty} \exp\left(-\frac{3}{8} \min\left\{\frac{t^2}{S^2}, \frac{t}{L}\right\}\right) dt \\
&\stackrel{(iii)}{\leq} \int_{\beta/2}^{\infty} \exp\left(-\frac{3}{8} \frac{t}{\max\{S/\sqrt{\log d}, L\}}\right) dt \\
&\lesssim \max\{S/\sqrt{\log d}, L\} \exp\left(-\frac{3}{16} \frac{\beta}{\max\{S/\sqrt{\log d}, L\}}\right) \\
&\stackrel{(iv)}{\lesssim} \beta \exp\left(-\frac{1}{16} C \log d\right) \ll \beta/d^2,
\end{aligned}$$

where (i) follows from (261), (ii) comes from (259), (iii) is a consequence of (262), and (iv) holds true when $C > 0$ is sufficiently large. Consequently,

$$\mathbb{E}\left[\max_{j,k \in [d]} \sum_{i \in [d]} R_{i,j,k}^2\right] \lesssim \beta \lesssim B\sqrt{\log d} + \sigma_{\text{mode}}.$$

Clearly, the same bound holds for $\mathbb{E}\left[\max_{i,k \in [d]} \sum_{j \in [d]} R_{i,j,k}^2\right]$ and $\mathbb{E}\left[\max_{i,j \in [d]} \sum_{k \in [d]} R_{i,j,k}^2\right]$.

Substitution into (257) yields

$$\mathbb{E}[\|\mathbf{R}\|] \lesssim B \log^3 d + \sigma_{\text{mode}} \log^{5/2} d. \quad (263)$$

Recognizing that the magnitudes of all entries of \mathbf{R} are bounded by B , we can invoke Talagrand's concentration inequality [Ver18, Theorem 5.2.16] for convex Lipschitz functions of independent bounded random variables to show that with probability $1 - O(d^{-10})$,

$$\left|\|\mathbf{R}\| - \mathbb{E}[\|\mathbf{R}\|]\right| \lesssim B\sqrt{\log d}$$

and, therefore,

$$\|\mathbf{R}\| \lesssim B \log^3 d + \sigma_{\text{mode}} \log^{5/2} d. \quad (264)$$

D.3.2 Proof of Corollary D.3

Now we apply Lemma D.2 to our concrete setting. We first look at $p^{-1}\mathcal{P}_{\Omega}(\mathbf{T}^*) - \mathbf{T}^*$ and treat it as \mathbf{R} in Lemma D.2. With the help of Lemma D.1, it is straightforward to compute that

$$\max_{i,j,k \in [d]} |T_{i,j,k}^* (p^{-1}\chi_{i,j,k} - 1)| \lesssim \frac{1}{p} \|\mathbf{A}^*\|_{\infty} \lesssim \frac{\sqrt{\mu r} \lambda_{\max}^*}{d^{3/2} p},$$

and

$$\max_{i,j \in [d]} \sum_{k \in [d]} \mathbb{E}[T_{i,j,k}^{*2} (p^{-1}\chi_{i,j,k} - 1)^2] \lesssim \frac{1}{p} \|\mathbf{A}^{*\top}\|_{2,\infty}^2 \lesssim \frac{\mu^2 r \lambda_{\max}^{*2}}{d^2 p}.$$

Clearly, $\max_{i,k \in [d]} \sum_{j \in [d]} \mathbb{E}[T_{i,j,k}^{*2} (p^{-1}\chi_{i,j,k} - 1)^2]$ and $\max_{j,k \in [d]} \sum_{i \in [d]} \mathbb{E}[T_{i,j,k}^{*2} (p^{-1}\chi_{i,j,k} - 1)^2]$ can be controlled in the same way. Substitution into (264) proves the claim (246).

We then turn to $\mathcal{P}_{\Omega}(\mathbf{E})$. Recognizing that the entries of \mathbf{E} might be unbounded, we invoke the following truncation trick to cope with this unboundedness issue. Specifically, define $\tilde{\mathbf{E}} = [\tilde{E}_{i,j,k}]_{1 \leq i,j,k \leq d}$ where

$$\tilde{E}_{i,j,k} := E_{i,j,k} \mathbf{1}\{|E_{i,j,k}| \leq c_1 \sigma \sqrt{\log d}\}, \quad 1 \leq i, j, k \leq d \quad (265)$$

for some sufficiently large constant $c_1 > 0$. Moreover, $\tilde{E}_{i,j,k}$ is zero-mean because we assume that the distribution of $E_{i,j,k}$ is symmetric about 0. Standard concentration inequalities reveal that: with probability exceeding $1 - O(d^{-10})$, one has $\mathbf{E} = \tilde{\mathbf{E}}$. Hence, it suffices to bound $\|\mathcal{P}_\Omega(\tilde{\mathbf{E}})\|$. Towards this end, simple calculation reveals that

$$\begin{aligned} B &= \max_{i,j,k \in [d]} |\tilde{E}_{i,j,k} \chi_{i,j,k}| \lesssim \|\tilde{\mathbf{E}}\|_\infty \lesssim \sigma \sqrt{\log d}, \\ \sigma_{\text{mode}}^2 &\leq \max_{i,j \in [d]} \sum_{k \in [d]} \mathbb{E}[E_{i,j,k}^2 \chi_{i,j,k}] + \max_{i,k \in [d]} \sum_{j \in [d]} \mathbb{E}[E_{i,j,k}^2 \chi_{i,j,k}] + \max_{j,k \in [d]} \sum_{i \in [d]} \mathbb{E}[E_{i,j,k}^2 \chi_{i,j,k}] \lesssim p\sigma^2 d. \end{aligned}$$

This together with (264) as well as the high-probability event $\mathbf{E} = \tilde{\mathbf{E}}$ completes the proof.

D.4 Proof of Lemma D.4

For notational simplicity, let us denote

$$\mathbf{X} := (p^{-1}\mathbf{T} - \mathbf{T}^*) \times_3 \mathbf{w}.$$

Observe that \mathbf{X} is a zero-mean random matrix in $\mathbb{R}^{d \times d}$ with independent entries

$$X_{i,j} = \sum_{k \in [d]} w_k \{T_{i,j,k}^* (p^{-1}\chi_{i,j,k} - 1) + p^{-1}E_{i,j,k} \chi_{i,j,k}\}, \quad (i,j) \in [d]^2.$$

We shall apply the truncated matrix Bernstein inequality to control the spectral norm of \mathbf{X} .

- First, it is straightforward to bound

$$\begin{aligned} V &:= \max \left\{ \max_{i \in [d]} \sum_{j \in [d]} \mathbb{E}[X_{i,j}^2], \max_{j \in [d]} \sum_{i \in [d]} \mathbb{E}[X_{i,j}^2] \right\} \\ &= \max \left\{ \max_{i \in [d]} \sum_{j,k \in [d]} p^{-1}w_k^2 (T_{i,j,k}^{*2} + \mathbb{E}[E_{i,j,k}^2]), \max_{j \in [d]} \sum_{i,k \in [d]} p^{-1}w_k^2 (T_{i,j,k}^{*2} + \mathbb{E}[E_{i,j,k}^2]) \right\} \\ &\leq \frac{1}{p} \left(\|\mathbf{w}\|_\infty^2 \|\mathbf{A}^*\|_{2,\infty}^2 + \|\mathbf{w}\|_2^2 \sigma^2 d \right). \end{aligned}$$

- Second, using the same truncation argument as in the proof of Lemma D.2 in Appendix D.3, we can assume $|E_{i,j,k}| \lesssim \sigma \sqrt{\log d}$ for all $1 \leq i, j, k \leq d$ (which holds with very high probability). The Bernstein inequality reveals that

$$\mathbb{P}\{|X_{i,j}| > t\} \leq 2 \exp\left(-\frac{3}{8} \min\left\{\frac{t^2}{S^2}, \frac{t}{L}\right\}\right), \quad t > 0$$

for each $(i,j) \in [d]^2$, where

$$\begin{aligned} L &:= \max_{k \in [d]} \{|w_k| |T_{i,j,k}^* (p^{-1}\chi_{i,j,k} - 1) + p^{-1}E_{i,j,k} \chi_{i,j,k}|\} \lesssim \frac{1}{p} \|\mathbf{w}\|_\infty \left(\|\mathbf{A}^*\|_\infty + \sigma \sqrt{\log d} \right); \\ S^2 &:= \mathbb{E}[X_{i,j}^2] \asymp \sum_{k \in [d]} p^{-1}w_k^2 (T_{i,j,k}^{*2} + \mathbb{E}[E_{i,j,k}^2]) \leq \frac{1}{p} \left(\|\mathbf{w}\|_\infty^2 \|\mathbf{A}^{*\top}\|_{2,\infty}^2 + \|\mathbf{w}\|_2^2 \sigma^2 \right). \end{aligned}$$

This implies that with probability exceeding $1 - O(d^{-20})$,

$$\begin{aligned} \max_{i,j \in [d]} |X_{i,j}| &\lesssim L \log d + S \sqrt{\log d} \\ &\lesssim \frac{\|\mathbf{w}\|_\infty \log d}{p} \left(\|\mathbf{A}^*\|_\infty + \sigma \sqrt{\log d} \right) + \sqrt{\frac{\log d}{p}} \left(\|\mathbf{w}\|_\infty \|\mathbf{A}^{*\top}\|_{2,\infty} + \|\mathbf{w}\|_2 \sigma \right). \end{aligned}$$

Therefore, if we choose

$$\beta := C \left\{ \frac{\|\mathbf{w}\|_\infty \log d}{p} \left(\|\mathbf{A}^*\|_\infty + \sigma \sqrt{\log d} \right) + \sqrt{\frac{\log d}{p}} \left(\|\mathbf{w}\|_\infty \|\mathbf{A}^{*\top}\|_{2,\infty} + \|\mathbf{w}\|_2 \sigma \right) \right\}$$

for some sufficiently large constant $C > 0$, then one has

$$\beta \geq \frac{C}{2} \left(L \log d + S \sqrt{\log d} \right).$$

- Third, it is easy to bound

$$\begin{aligned} \mathbb{E}[|X_{i,j}| \mathbf{1}\{|X_{i,j}| \geq \beta\}] &\leq \beta \cdot \mathbb{P}\{|X_{i,j}| \geq \beta\} + \int_\beta^\infty \mathbb{P}\{|X_{i,j}| \geq t\} dt \\ &\leq \beta \cdot O(d^{-20}) + \int_\beta^\infty \mathbb{P}\{|X_{i,j}| \geq t\} dt. \end{aligned}$$

In view of our choice of β , we know that $\min\{t^2/S^2, t/L\} \geq t/\max\{S/\sqrt{\log d}, L\}$ for any $t \geq \beta$. As a result, for d sufficiently large, we have

$$\begin{aligned} \int_\beta^\infty \mathbb{P}\{|X_{i,j}| \geq t\} dt &\leq 2 \int_\beta^\infty \exp\left(-\frac{3}{8} \min\left\{\frac{t^2}{S^2}, \frac{t}{L}\right\}\right) dt \\ &\leq 2 \int_\beta^\infty \exp\left(-\frac{3}{8} \frac{t}{\max\{S/\sqrt{\log d}, L\}}\right) dt \\ &\lesssim \max\{S/\sqrt{\log d}, L\} \exp\left(-\frac{3}{8} \frac{\beta}{\max\{S/\sqrt{\log d}, L\}}\right) \\ &\lesssim \max\{S/\sqrt{\log d}, L\} \exp\left(-\frac{3}{8} C \log d\right) \ll \frac{\beta}{d^2}. \end{aligned}$$

Consequently, we have established that

$$q := \sum_{i,j} \mathbb{E}[|X_{i,j}| \mathbf{1}\{|X_{i,j}| \geq \beta\}] \ll \beta.$$

Invoke the matrix Bernstein inequality to demonstrate that with probability $1 - O(d^{-10})$,

$$\begin{aligned} \|\mathbf{X}\| &\lesssim q + \beta \log d + \sqrt{V \log d} \asymp \beta \log d + \sqrt{V \log d} \\ &\lesssim \frac{\|\mathbf{w}\|_\infty \log^2 d}{p} \left(\|\mathbf{A}^*\|_\infty + \sigma \sqrt{\log d} \right) + \frac{\log^{3/2} d}{\sqrt{p}} \left(\|\mathbf{w}\|_\infty \|\mathbf{A}^{*\top}\|_{2,\infty} + \sigma \|\mathbf{w}\|_2 \right) \\ &\quad + \sqrt{\frac{\log d}{p}} \left(\|\mathbf{w}\|_\infty \|\mathbf{A}^*\|_{2,\infty} + \|\mathbf{w}\|_2 \sigma \sqrt{d} \right) \\ &\asymp \|\mathbf{w}\|_\infty \left\{ \frac{\|\mathbf{A}^*\|_\infty \log^2 d}{p} + \frac{\|\mathbf{A}^{*\top}\|_{2,\infty} \log^{3/2} d}{\sqrt{p}} + \frac{\|\mathbf{A}^*\|_{2,\infty} \sqrt{\log d}}{\sqrt{p}} + \frac{\sigma \log^{5/2} d}{p} \right\} + \|\mathbf{w}\|_2 \sigma \sqrt{\frac{d \log d}{p}} \\ &\stackrel{(i)}{\lesssim} \|\mathbf{w}\|_\infty \left\{ \frac{\sqrt{\mu r} \log^2 d}{d^{3/2} p} + \sqrt{\frac{\mu^2 r \log^3 d}{d^2 p}} + \sqrt{\frac{\mu r \log d}{d p}} \right\} \lambda_{\max}^* + \|\mathbf{w}\|_\infty \frac{\sigma \log^{5/2} d}{p} + \|\mathbf{w}\|_2 \sigma \sqrt{\frac{d \log d}{p}} \\ &\stackrel{(ii)}{\lesssim} \|\mathbf{w}\|_\infty \sqrt{\frac{\mu r \log d}{d p}} \lambda_{\max}^* + \|\mathbf{w}\|_\infty \frac{\sigma \log^{5/2} d}{p} + \|\mathbf{w}\|_2 \sigma \sqrt{\frac{d \log d}{p}}, \end{aligned}$$

where (i) is due to Lemma D.1, and (ii) follows as long as $p \gtrsim d^{-2} \log^3 d$ and $\mu \log^2 d \lesssim d$.

D.5 Proof of Lemma D.5

Recall that for a standard Gaussian random variable $Z \sim \mathcal{N}(0, 1)$, one has

$$\frac{1}{5t\sqrt{2\pi}} \exp(-t^2/2) \leq \left(\frac{1}{t} - \frac{1}{t^3}\right) \frac{1}{\sqrt{2\pi}} \exp(-t^2/2) \leq \mathbb{P}\{Z \geq t\} \leq \frac{1}{t} \frac{1}{\sqrt{2\pi}} \exp(-t^2/2) \quad (266)$$

for all $t > \sqrt{5/4}$. Observing that $\kappa\sqrt{2\log r} + \Delta \geq \sqrt{5/4}$ since $\kappa \geq 1$ and $r \geq 2$, we can invoke the above tail bound to deduce that

$$\begin{aligned} \mathbb{P}\left\{X_{1,j} \geq \kappa\sqrt{2\log r} + \Delta\right\} &> \frac{1}{5\sqrt{2\pi}(\kappa\sqrt{2\log r} + \Delta)} \exp\left(-(\kappa\sqrt{2\log r} + \Delta)^2/2\right) \\ &\geq \frac{1}{5\sqrt{2\pi}(\kappa\sqrt{2\log r} + \Delta)r^{2\kappa^2} \exp(\Delta^2)}, \end{aligned} \quad (267)$$

where we use the elementary inequality $(\kappa\sqrt{2\log r} + \Delta)^2 \leq 4\kappa^2 \log r + 2\Delta^2$. In addition, it follows from the union bound that

$$\begin{aligned} \mathbb{P}\left\{\max_{1 < i \leq r} |X_{i,j}| < \sqrt{2\log r}\right\} &\geq 1 - r \mathbb{P}\left\{|X_{i,j}| > \sqrt{2\log r}\right\} \geq 1 - r \left\{\frac{1}{2\sqrt{\pi \log r}} \exp(-\log r)\right\} \\ &\geq 1 - \frac{1}{\sqrt{\pi \log r}} \geq 1 - \frac{1}{\sqrt{\pi \log 2}}. \end{aligned} \quad (268)$$

To prove the claim, it is sufficient to choose L such that

$$\mathbb{P}\left\{\forall j : X_{1,j} < \kappa\sqrt{2\log r} + \Delta \text{ or } \max_{1 < i \leq r} |X_{i,j}| \geq \sqrt{2\log r}\right\} \leq \delta,$$

or equivalently,

$$\left(1 - \mathbb{P}\{X_{1,j} \geq \kappa\sqrt{2\log r} + \Delta\} \mathbb{P}\{\max_{1 < i \leq r} |X_{i,j}| < \sqrt{2\log r}\}\right)^L \leq \delta. \quad (269)$$

Note that $\log(1-x) \leq -1/(2x)$ for $0 < x < 1/4$. In view of (267) and (268), one can verify that the above inequality (269) as long as

$$L \geq C(\kappa\sqrt{\log r} + \Delta)r^{2\kappa^2} \exp(\Delta^2) \log \frac{1}{\delta}, \quad (270)$$

where $C > 0$ is some universal constant.

To prove the second claim, recall the definitions that

$$\Delta_j := X_{1,j} - \max_{1 < i \leq r} \kappa X_{i,j}, \quad 1 \leq j \leq L.$$

and that $\Delta_{(1)} \geq \Delta_{(2)} \geq \dots \geq \Delta_{(L)}$ denote $\{\Delta_j\}_{j=1}^L$ in descending order. For any $\epsilon > 0$, one has

$$\begin{aligned} \mathbb{P}\{\Delta_{(1)} - \Delta_{(2)} < \epsilon\} &= \sum_{1 \leq j \leq L} \mathbb{P}\{\Delta_j - \max_{k:k \neq j} \Delta_k < \epsilon \mid \Delta_j = \Delta_{(1)}\} \mathbb{P}\{\Delta_j = \Delta_{(1)}\} \\ &= \mathbb{P}\{\Delta_1 - \max_{k:k \neq 1} \Delta_k > \epsilon \mid \Delta_1 = \Delta_{(1)}\}, \end{aligned}$$

where the last line holds because the distribution of $\Delta_j - \max_{k:k \neq j} \Delta_k$ conditional on $\Delta_j = \max_{1 \leq k \leq r} \Delta_k$ is identical for all $1 \leq j \leq L$. In addition, it is straightforward to see that

$$\begin{aligned} \Delta_1 = \Delta_{(1)} &\iff \Delta_1 \geq \max_{k:k \neq 1} \Delta_k \\ &\iff X_{1,1} \geq \max_{1 < i \leq r} \kappa X_{i,1} + \max_{k:k \neq 1} \Delta_k =: Y_{1,1}. \end{aligned}$$

Hence, we have

$$\mathbb{P}\{\Delta_{(1)} - \Delta_{(2)} < \epsilon\} = \mathbb{P}\{\Delta_1 - \max_{k:k \neq 1} \Delta_k < \epsilon \mid \Delta_1 = \Delta_{(1)}\} = \mathbb{P}\{X_{1,1} - Y_{1,1} < \epsilon \mid X_{1,1} \geq Y_{1,1}\}.$$

Next, observe that $X_{1,1}$ is independent of $Y_{1,1}$, and hence we have

$$\begin{aligned} \mathbb{P}\{X_{1,1} - Y_{1,1} < \epsilon \mid X_{1,1} \geq Y_{1,1}, Y_{1,1} = x\} &= \mathbb{P}\{X_{1,1} < x + \epsilon \mid X_{1,1} \geq y, Y_{1,1} = x\} \\ &= \mathbb{P}\{X_{1,1} < x + \epsilon \mid X_{1,1} \geq x\} = \frac{\mathbb{P}\{x \leq X_{1,1} < x + \epsilon\}}{\mathbb{P}\{X_{1,1} \geq x\}} \end{aligned}$$

for any $x \geq 0$. In order to study this function, we define $f_\epsilon(x) := \frac{\mathbb{P}\{x \leq Z \leq x + \epsilon\}}{\mathbb{P}\{Z \geq x\}}$ with $Z \sim \mathcal{N}(0, 1)$. Taking the derivative of $f_\epsilon(\cdot)$ w.r.t. x gives: for any $\epsilon > 0$,

$$\begin{aligned} f'_\epsilon(x) &= \frac{\{\exp(-(x + \epsilon)^2/2) - \exp(-x^2/2)\} \mathbb{P}\{Z \geq x\} + \exp(-x^2/2) \mathbb{P}\{x \leq Z \leq x + \epsilon\}}{\sqrt{2\pi} (\mathbb{P}\{Z \geq x\})^2} \\ &= \frac{\exp(-(x + \epsilon)^2/2) \mathbb{P}\{Z \geq x\} - \exp(-\frac{x^2}{2}) \mathbb{P}\{Z \geq x + \epsilon\}}{\sqrt{2\pi} (\mathbb{P}\{N \geq x\})^2} \\ &= \frac{\exp(x^2/2) \mathbb{P}\{Z \geq x\} - \exp((x + \epsilon)^2/2) \mathbb{P}\{Z \geq x + \epsilon\}}{\sqrt{2\pi} \exp((x + \epsilon)^2/2) \exp(x^2/2) (\mathbb{P}\{Z \geq x\})^2} \\ &= \frac{\int_0^\infty \left(\exp\left(-\frac{t^2 + 2tx}{2}\right) - \exp\left(-\frac{t^2 + 2t(x + \epsilon)}{2}\right) \right) dt}{2\pi \exp((x + \epsilon)^2/2) \exp(x^2/2) (\mathbb{P}\{Z \geq x\})^2} > 0. \end{aligned}$$

In other words, $f_\epsilon(x)$ is monotonically increasing in x for any given $\epsilon > 0$. Therefore, for any $0 \leq x < B$ for some sufficiently large $B > \sqrt{5/4}$ (to be specified later), the above bounds taken together give

$$\mathbb{P}\{X_{1,1} - Y_{1,1} < \epsilon \mid X_{1,1} \geq Y_{1,1}, Y_{1,1} = x\} \stackrel{(i)}{\leq} \frac{\mathbb{P}\{B \leq X_{1,1} < B + \epsilon\}}{\mathbb{P}\{X_{1,1} \geq B\}} \stackrel{(ii)}{\leq} \frac{\epsilon \exp(-B^2/2)}{\frac{1}{5B} \exp(-B^2/2)} = 5\epsilon B,$$

where (i) arises from the monotonicity of $f_\epsilon(\cdot)$, and (ii) relies on (266). By taking $\epsilon = \delta/(5B)$, we obtain

$$\mathbb{P}\{X_{1,1} - Y_{1,1} < \epsilon \mid X_{1,1} \geq Y_{1,1}, Y_{1,1} = x\} \leq \delta$$

for any $0 \leq x \leq B$. Recall that $Y_{1,1} = \max_{1 < i \leq r} \kappa X_{i,1} + \max_{k:k \neq 1} \Delta_k$. By standard Gaussian concentration inequalities, with probability at least $1 - \delta$ one has

$$Y_{1,1} \lesssim \kappa \sqrt{\log r} + \sqrt{\log L} + \sqrt{\log(1/\delta)} \asymp \sqrt{\log L} + \sqrt{\log(1/\delta)},$$

where the last step arises from the lower bound on L in (270). If we choose $B = C(\sqrt{\log L} + \sqrt{\log(1/\delta)})$ for some sufficiently large universal constant $C > 0$, then this immediately implies that

$$\begin{aligned} \mathbb{P}\{\Delta_{(1)} - \Delta_{(2)} < \epsilon\} &= \mathbb{P}\{X_{1,1} - Y_{1,1} < \epsilon \mid X_{1,1} \geq Y_{1,1}\} \\ &\leq \mathbb{P}\{Y_{1,1} > B\} + \mathbb{P}\{X_{1,1} - Y_{1,1} < \epsilon \mid X_{1,1} \geq Y_{1,1}, Y_{1,1} = B\} \leq 2\delta. \end{aligned}$$

We have therefore concluded the proof.

D.6 Proof of Lemma D.6

To begin with, it is self-evident that

$$\mathcal{P}_{\mathbf{V}}(\mathbf{u}_0) = \mathbf{V}\mathbf{V}^\top \mathbf{u}_0 = \mathbf{U}\mathbf{U}^\top \mathbf{u}_0 + (\mathbf{V}\mathbf{V}^\top - \mathbf{U}\mathbf{U}^\top) \mathbf{u}_0 = \mathbf{u}_0 + (\mathbf{V}\mathbf{V}^\top - \mathbf{U}\mathbf{U}^\top) \mathbf{u}_0,$$

where the last identity follows since \mathbf{u}_0 is assumed to lie within $\text{span}(\mathbf{U})$. As a result,

$$\begin{aligned} \mathcal{P}_{\mathbf{V}^\perp}(\mathbf{u}_0) &= \mathbf{u}_0 - \mathcal{P}_{\mathbf{V}}(\mathbf{u}_0) = -(\mathbf{V}\mathbf{V}^\top - \mathbf{U}\mathbf{U}^\top) \mathbf{u}_0 \\ \implies \|\mathcal{P}_{\mathbf{V}^\perp}(\mathbf{u}_0)\|_2 &\leq \|\mathbf{V}\mathbf{V}^\top - \mathbf{U}\mathbf{U}^\top\| \cdot \|\mathbf{u}_0\|_2 \leq \delta. \end{aligned}$$

The Pythagorean theorem then gives $\|\mathcal{P}_{\mathbf{V}}(\mathbf{u}_0)\|_2 = \sqrt{\|\mathbf{u}_0\|_2^2 - \|\mathcal{P}_{\mathbf{V}^\perp}(\mathbf{u}_0)\|_2^2} \geq \sqrt{1 - \delta^2}$.

D.7 Proof of Lemma D.7

By virtue of Lemma D.1, we can compute

$$\begin{aligned} \sum_{i,j \in [d]} \mathbb{E} [T_{i,j,k}^{*2} (p^{-1} \chi_{i,j,k} - 1)^2] &\leq \frac{1}{p} \sum_{i,j \in [d]} T_{i,j,k}^{*2} \leq \frac{1}{p} \|\mathbf{A}^*\|_{2,\infty}^2 \lesssim \frac{\mu r \lambda_{\max}^{*2}}{dp} =: M; \\ \left| T_{i,j,k}^{*2} (p^{-1} \chi_{i,j,k} - 1)^2 \right| &\leq \frac{1}{p^2} \|\mathbf{T}^*\|_{\infty}^2 = \frac{1}{p^2} \|\mathbf{A}^*\|_{\infty}^2 \lesssim \frac{\mu r \lambda_{\max}^{*2}}{d^3 p^2} =: L; \\ \sum_{i,j \in [d]} \text{Var} \left(T_{i,j,k}^{*2} (p^{-1} \chi_{i,j,k} - 1)^2 \right) &\lesssim \frac{1}{p^3} \sum_{i,j \in [d]} T_{i,j,k}^{*4} \leq \frac{1}{p^3} \|\mathbf{A}^*\|_{\infty}^2 \|\mathbf{A}^*\|_{2,\infty}^2 \lesssim \frac{\mu^2 r^2 \lambda_{\max}^{*4}}{d^4 p^3} =: V. \end{aligned}$$

Invoke the Bernstein inequality to show that: with probability exceeding $1 - O(d^{-20})$,

$$\begin{aligned} \sum_{i,j \in [d]} T_{i,j,k}^{*2} (p^{-1} \chi_{i,j,k} - 1)^2 &\lesssim M + L \log d + \sqrt{V \log d} \\ &\lesssim \frac{\mu r \lambda_{\max}^{*2}}{dp} + \frac{\mu r \lambda_{\max}^{*2} \log d}{d^3 p^2} + \sqrt{\frac{\mu^2 r^2 \lambda_{\max}^{*4} \log d}{d^4 p^3}} \\ &\asymp \frac{\mu r \lambda_{\max}^{*2}}{dp}, \end{aligned}$$

where the last line holds with the proviso that $p \gtrsim d^{-2} \log d$.

D.8 Proof of Lemma D.8

Since the $E_{i,j,k}$'s are independent sub-Gaussian random variables with variance at most σ^2 , one has

$$\begin{aligned} \sum_{i,j \in [d]} \mathbb{E} [(E_{i,j,k} \chi_{i,j,k})^2] &\lesssim \sigma^2 d^2 p =: M; \\ \|(E_{i,j,k} \chi_{i,j,k})^2\|_{\psi_1} &\lesssim \sigma^2 =: L; \\ \sum_{i,j \in [d]} \text{Var} [(E_{i,j,k} \chi_{i,j,k})^2] &\lesssim \sigma^4 d^2 p =: V. \end{aligned}$$

Here, $\|\cdot\|_{\psi_1}$ denotes the sub-exponential norm [Ver10]. Taken together with the Bernstein inequality, these yield that with probability exceeding $1 - O(d^{-20})$,

$$\sum_{i,j \in [d]} (E_{i,j,k} \chi_{i,j,k})^2 \lesssim M + \sqrt{V \log d} + L \log^2 d \lesssim \sigma^2 d^2 p + \sqrt{\sigma^4 d^2 p \log d} + \sigma^2 \log^2 d \asymp \sigma^2 d^2 p,$$

provided that $p \gtrsim d^{-2} \log^2 d$.

D.9 Proof of Lemma D.9

Fix an arbitrary $1 \leq i \leq d$. We first define a sequence of independent zero-mean random variables $\{X_j\}_{j \in [d]}$ as follows

$$X_j := \sum_{k \in [d]} T_{i,j,k}^* w_k (p^{-1} \chi_{i,j,k} - 1).$$

One can easily show that

$$\begin{aligned} \max_{k \in [d]} \left| T_{i,j,k}^* w_k (p^{-1} \chi_{i,j,k} - 1) \right| &\leq \frac{1}{p} \|\mathbf{A}^*\|_{\infty} \|\mathbf{w}\|_{\infty} =: L, \\ \mathbb{E}[X_j^2] &= \sum_{k \in [d]} T_{i,j,k}^{*2} w_k^2 \mathbb{E}[(p^{-1} \chi_{i,j,k} - 1)^2] \leq \frac{1}{p} \|\mathbf{w}\|_{\infty}^2 \sum_{k \in [d]} T_{i,j,k}^{*2} = \frac{1}{p} \|\mathbf{w}\|_{\infty}^2 \|\mathbf{A}^{*\top}\|_{2,\infty}^2 =: V. \end{aligned}$$

The Bernstein inequality indicates that: with probability at least $1 - O(d^{-20})$,

$$|X_j| \lesssim L \log d + \sqrt{V \log d} \lesssim \frac{1}{p} \|\mathbf{A}^*\|_\infty \|\mathbf{w}\|_\infty \log d + \sqrt{\frac{\log d}{p}} \|\mathbf{A}^{*\top}\|_{2,\infty} \|\mathbf{w}\|_\infty := L_j. \quad (271)$$

Moreover, we can also bound the variance of X_j as follows

$$\begin{aligned} \text{Var}(X_j^2) &\leq \mathbb{E}[X_j^4] \lesssim \frac{1}{p^3} \sum_{k \in [d]} T_{i,j,k}^{*4} w_k^4 + \frac{1}{p^2} \sum_{k_1 \neq k_2} T_{i,j,k_1}^{*2} T_{i,j,k_2}^{*2} w_{k_1}^2 w_{k_2}^2 \\ &\lesssim \frac{1}{p^3} \|\mathbf{w}\|_\infty^4 \|\mathbf{A}^*\|_\infty^2 \sum_{k \in [d]} T_{i,j,k}^{*2} + \frac{1}{p^2} \|\mathbf{w}\|_\infty^4 \|\mathbf{A}^{*\top}\|_{2,\infty}^2 \sum_{k \in [d]} T_{i,j,k}^{*2}. \end{aligned}$$

Given that X_j might be overly large in some rare case, we introduce a sequence $\{Y_j\}$, where we denote by Y_j the truncated version of X_j as follows

$$Y_j := X_j \mathbf{1}\{|X_j| \lesssim L_j\}.$$

We have learn from (271) and the union bound that with probability at least $1 - O(d^{-15})$, one has $Y_j = X_j$ for all $j \in [d]$.

Using the above bounds on the X_j 's, one observes that $\{Y_j\}_{j \in [d]}$ is a sequence of independent random variables satisfying

$$\begin{aligned} \sum_{j \in [d]} \mathbb{E}[Y_j^2] &\leq \sum_{j \in [d]} \mathbb{E}[X_j^2] \leq \frac{1}{p} \|\mathbf{A}^*\|_{2,\infty}^2 \|\mathbf{w}\|_\infty^2; \\ \max_{j \in [d]} Y_j^2 &\leq \max_{j \in [d]} L_j^2 \lesssim \frac{1}{p^2} \|\mathbf{A}^*\|_\infty^2 \|\mathbf{w}\|_\infty^2 \log d + \frac{\log d}{p} \|\mathbf{A}^{*\top}\|_{2,\infty}^2 \|\mathbf{w}\|_\infty^2; \\ \sum_{j \in [d]} \text{Var}(Y_j^2) &\leq \sum_{j \in [d]} \mathbb{E}[X_j^4] \lesssim \frac{1}{p^3} \|\mathbf{A}^*\|_\infty^2 \|\mathbf{A}^*\|_{2,\infty}^2 \|\mathbf{w}\|_\infty^4 + \frac{1}{p^2} \|\mathbf{A}^{*\top}\|_{2,\infty}^2 \|\mathbf{A}^*\|_{2,\infty}^2 \|\mathbf{w}\|_\infty^4. \end{aligned}$$

We can apply the Bernstein inequality to conclude that: with probability greater than $1 - O(d^{-15})$,

$$\begin{aligned} \sum_{j \in [d]} Y_j^2 &\lesssim \sum_{j \in [d]} \mathbb{E}[Y_j^2] + \max_{j \in [d]} Y_j^2 \log d + \sqrt{\sum_{j \in [d]} \text{Var}(Y_j^2) \log d} \\ &\lesssim \frac{1}{p} \|\mathbf{A}^*\|_{2,\infty}^2 \|\mathbf{w}\|_\infty^2 + \frac{\log^2 d}{p^2} \|\mathbf{A}^*\|_\infty^2 \|\mathbf{w}\|_\infty^2 + \frac{\log^2 d}{p} \|\mathbf{A}^{*\top}\|_{2,\infty}^2 \|\mathbf{w}\|_\infty^2 \\ &\quad + \sqrt{\frac{\log d}{p^3} \|\mathbf{A}^*\|_\infty^2 \|\mathbf{A}^*\|_{2,\infty}^2 \|\mathbf{w}\|_\infty^4 + \frac{\log d}{p^2} \|\mathbf{w}\|_\infty^4 \|\mathbf{A}^{*\top}\|_{2,\infty}^2 \|\mathbf{A}^*\|_{2,\infty}^2} \\ &\stackrel{(i)}{\lesssim} \frac{1}{p} \|\mathbf{A}^*\|_{2,\infty}^2 \|\mathbf{w}\|_\infty^2 + \frac{\log^2 d}{p^2} \|\mathbf{A}^*\|_\infty^2 \|\mathbf{w}\|_\infty^2 + \frac{\log^2 d}{p} \|\mathbf{A}^{*\top}\|_{2,\infty}^2 \|\mathbf{w}\|_\infty^2 \\ &\stackrel{(ii)}{\lesssim} \left(\frac{\mu r \lambda_{\max}^{*2}}{dp} + \frac{\mu r \lambda_{\max}^{*2} \log^2 d}{d^3 p^2} + \frac{\mu^2 r \lambda_{\max}^{*2} \log^2 d}{d^2 p} \right) \|\mathbf{w}\|_\infty^2 \\ &\asymp \frac{\mu r \lambda_{\max}^{*2}}{dp} \|\mathbf{w}\|_\infty^2, \end{aligned}$$

where (i) is due to the AM-GM inequality, (ii) makes use of Lemma D.1, and the last line follows under the conditions $p \gtrsim d^{-2} \log^2 d$ and $\mu \log^2 d \lesssim d$. This together with the high-probability fact $Y_j = X_j$ ($\forall j \in [d]$) concludes the proof.

D.10 Proof of Lemma D.10

Fix any $1 \leq i \leq d$. To begin with, define

$$Z_j := \sum_{k \in [d]} w_k E_{i,j,k} \chi_{i,j,k},$$

which is a zero-mean random variable. In order to bound Z_j , one observes that

$$\begin{aligned} \|w_k E_{i,j,k} \chi_{i,j,k}\|_{\psi_1} &\lesssim \sigma \|\mathbf{w}\|_\infty =: L; \\ \sum_{k \in [d]} \text{Var}(w_k E_{i,j,k} \chi_{i,j,k}) &\leq \sigma^2 p \sum_{k \in [d]} w_k^2 = \sigma^2 p \|\mathbf{w}\|_2^2 =: V, \end{aligned}$$

where $\|\cdot\|_{\psi_1}$ denotes the sub-exponential norm. Apply the Bernstein inequality for the sum of sub-exponential random variables to obtain

$$|Z_j| \lesssim \sqrt{V \log d} + L \log^2 d \lesssim \sigma \|\mathbf{w}\|_2 \sqrt{p \log d} + \sigma \|\mathbf{w}\|_\infty \log^2 d =: L_j \quad (272)$$

with probability exceeding $1 - O(d^{-20})$. Further, given that Z_j is not necessarily bounded, we introduce a sequence of truncated random variables as follows

$$Y_j := Z_j \mathbf{1}\{|Z_j| \lesssim L_j\}. \quad (273)$$

According to the above bound, one has $Y_j = Z_j$ ($\forall j$) with probability at least $1 - O(d^{-19})$.

We then turn attention to bounding $\sum_{j \in [d]} Y_j^2$. To this end, observe that

$$\begin{aligned} \sum_{j \in [d]} \mathbb{E}[Y_j^2] &\leq \sum_{j \in [d]} \mathbb{E}\left[\left(\sum_{k \in [d]} w_k E_{i,j,k} \chi_{i,j,k}\right)^2\right] \lesssim \sigma^2 p \sum_{j \in [d]} \sum_{k \in [d]} w_k^2 \\ &= \sigma^2 p d \|\mathbf{w}\|_2^2 =: M_0. \end{aligned}$$

Additionally,

$$\begin{aligned} \sum_{j \in [d]} \mathbb{E}[Y_j^4] &\leq \sum_{j \in [d]} \mathbb{E}\left[\left(\sum_{k \in [d]} w_k E_{i,j,k} \chi_{i,j,k}\right)^4\right] \\ &\leq \sum_{j \in [d]} \sum_{k \in [d]} \mathbb{E}\left[w_k^4 E_{i,j,k}^4 \chi_{i,j,k}^4\right] + \sum_{j \in [d]} \sum_{k_1 \neq k_2} \mathbb{E}\left[w_{k_1}^2 w_{k_2}^2 E_{i,j,k_1}^2 E_{i,j,k_2}^2 \chi_{i,j,k_1}^2 \chi_{i,j,k_2}^2\right] \\ &\lesssim \sigma^4 p \sum_{j \in [d]} \sum_{k \in [d]} w_k^4 + \sigma^4 p^2 \sum_{j \in [d]} \sum_{1 \leq k_1 \neq k_2 \leq d} w_{k_1}^2 w_{k_2}^2 \\ &\lesssim \sigma^4 p d \|\mathbf{w}\|_2^2 \|\mathbf{w}\|_\infty^2 + \sigma^4 p^2 d \|\mathbf{w}\|_2^4 =: V_0. \end{aligned}$$

Invoke the Bernstein inequality to arrive at: with probability at least $1 - O(d^{-20})$,

$$\begin{aligned} \sum_{j \in [d]} Y_j^2 &\lesssim M_0 + \sqrt{V_0 \log d} + \max_{j \in [d]} L_j^2 \log d \\ &\lesssim \sigma^2 p d \|\mathbf{w}\|_2^2 + \sqrt{\sigma^4 p d \|\mathbf{w}\|_2^2 \left(p \|\mathbf{w}\|_2^2 + \|\mathbf{w}\|_\infty^2\right) \log d} + \left(\sigma^2 p \|\mathbf{w}\|_2^2 \log^2 d + \sigma^2 \|\mathbf{w}\|_\infty^2 \log^5 d\right) \\ &\asymp \sigma^2 p d \|\mathbf{w}\|_2^2 + \sigma^2 \|\mathbf{w}\|_\infty^2 \log^5 d. \end{aligned}$$

This together with the high-probability fact $Y_j = Z_j$ ($\forall j$) completes the proof.

E Extension to asymmetric tensors

Thus far, we have focused on the case where the tensor of interest is symmetric. In this section, we discuss how to generalize our algorithm and analysis to accommodate asymmetric tensors.

E.1 Problem settings

Suppose that the unknown tensor $\mathbf{T}^* \in \mathbb{R}^{d_1 \times d_2 \times d_3}$ is a rank- r tensor with CP decomposition

$$\mathbf{T}^* = \sum_{i=1}^r \mathbf{u}_i^* \otimes \mathbf{v}_i^* \otimes \mathbf{w}_i^*, \quad (274)$$

where $\mathbf{u}_i^* \in \mathbb{R}^{d_1}$, $\mathbf{v}_i^* \in \mathbb{R}^{d_2}$, $\mathbf{w}_i^* \in \mathbb{R}^{d_3}$ represent the tensor factors of interest. Apparently, there is an unavoidable global scaling ambiguity issue (for instance, multiplying \mathbf{u}_i^* by a constant c and multiplying \mathbf{v}_i^* by $1/c$ accordingly result in the same tensor). Without loss of generality, we shall assume throughout that

$$\|\mathbf{u}_i^*\|_2 = \|\mathbf{v}_i^*\|_2 = \|\mathbf{w}_i^*\|_2, \quad 1 \leq i \leq r. \quad (275)$$

In addition, we assume that each entry (j, k, l) is included in the sampling set Ω independently with probability p , and that each observed entry $T_{j,k,l}^*$ is corrupted by an independent zero-mean sub-Gaussian noise $E_{j,k,l}$ (cf. Assumption 2.3). Our goal is to (1) estimate $\{\mathbf{u}_i^*, \mathbf{v}_i^*, \mathbf{w}_i^*\}_{i=1}^r$ faithfully, modulo global permutation and global signs, and (2) estimate \mathbf{T}^* in a reliable manner.

E.2 Algorithms

We now move on to present an extension of our nonconvex algorithm to handle the noisy scenario.

First of all, setting

$$\mathbf{U} := [\mathbf{u}_1, \dots, \mathbf{u}_r] \in \mathbb{R}^{d_1 \times r}, \quad \mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_r] \in \mathbb{R}^{d_2 \times r} \quad \text{and} \quad \mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_r] \in \mathbb{R}^{d_3 \times r},$$

we can define the following regularized squared loss function

$$g(\mathbf{U}, \mathbf{V}, \mathbf{W}) := \frac{1}{6p} \left\| \mathcal{P}_\Omega \left(\sum_{i=1}^r \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i - \mathbf{T} \right) \right\|_{\mathbb{F}}^2 + \text{reg}(\mathbf{U}, \mathbf{V}, \mathbf{W}), \quad (276)$$

where the regularization term $\text{reg}(\mathbf{U}, \mathbf{V}, \mathbf{W})$ is given by

$$\text{reg}(\mathbf{U}, \mathbf{V}, \mathbf{W}) := \frac{1}{24} \sum_{i=1}^r \alpha_i \left\{ (\|\mathbf{u}_i\|_2^2 - \|\mathbf{v}_i\|_2^2)^2 + (\|\mathbf{u}_i\|_2^2 - \|\mathbf{w}_i\|_2^2)^2 + (\|\mathbf{v}_i\|_2^2 - \|\mathbf{w}_i\|_2^2)^2 \right\} \quad (277)$$

for some positive regularization parameters $\{\alpha_i\}_{i=1}^r$ to be specified momentarily. In contrast to the symmetric case, the addition term $\text{reg}(\mathbf{U}, \mathbf{V}, \mathbf{W})$ is included to help ensure that the sizes of \mathbf{U} , \mathbf{V} and \mathbf{W} stay close — an algorithmic trick that has proved useful in other problems like nonconvex rectangular matrix recovery [TBS⁺16, ZL16, CLL19].

We are now ready to present our nonconvex algorithm that accommodates the case with asymmetric tensors. As before, the proposed algorithm is initialized by a spectral method, followed by gradient descent designed to minimize the regularized loss function (276). The precise procedure is described in Algorithm 7 (which invokes Algorithms 8-9).

Before proceeding, we find it helpful to record closed-form expressions for the gradients, which are a crucial part when implementing the nonconvex gradient descent algorithm. Specifically, the gradients of $g(\mathbf{U}, \mathbf{V}, \mathbf{W})$ can be computed as follows

$$\nabla_{\mathbf{u}_i} g(\mathbf{U}, \mathbf{V}, \mathbf{W}) = \frac{1}{3p} \mathcal{P}_\Omega \left(\sum_{i=1}^r \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i - \mathbf{T} \right) \times_2 \mathbf{v}_i \times_3 \mathbf{w}_i + \frac{1}{6} \alpha_i (2 \|\mathbf{u}_i\|_2^2 - \|\mathbf{v}_i\|_2^2 - \|\mathbf{w}_i\|_2^2) \mathbf{u}_i, \quad (278a)$$

$$\nabla_{\mathbf{v}_i} g(\mathbf{U}, \mathbf{V}, \mathbf{W}) = \frac{1}{3p} \mathcal{P}_\Omega \left(\sum_{i=1}^r \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i - \mathbf{T} \right) \times_1 \mathbf{u}_i \times_3 \mathbf{w}_i + \frac{1}{6} \alpha_i (2 \|\mathbf{v}_i\|_2^2 - \|\mathbf{u}_i\|_2^2 - \|\mathbf{w}_i\|_2^2) \mathbf{v}_i, \quad (278b)$$

$$\nabla_{\mathbf{w}_i} g(\mathbf{U}, \mathbf{V}, \mathbf{W}) = \frac{1}{3p} \mathcal{P}_\Omega \left(\sum_{i=1}^r \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i - \mathbf{T} \right) \times_1 \mathbf{u}_i \times_2 \mathbf{v}_i + \frac{1}{6} \alpha_i (2 \|\mathbf{w}_i\|_2^2 - \|\mathbf{u}_i\|_2^2 - \|\mathbf{v}_i\|_2^2) \mathbf{w}_i \quad (278c)$$

for each $1 \leq i \leq r$, where \times_1 , \times_2 and \times_3 have been defined in Section 2.4.

E.3 Numerical experiments

In order to validate the effectiveness of the proposed algorithm, we conduct a series of numerical experiments.

Algorithm 7 Gradient descent for nonconvex tensor completion (asymmetric case)

- 1: Generate initial estimates $\mathbf{U}^0 \in \mathbb{R}^{d_1 \times r}$, $\mathbf{V}^0 \in \mathbb{R}^{d_2 \times r}$, $\mathbf{W}^0 \in \mathbb{R}^{d_3 \times r}$ via Algorithm 8.
- 2: **for** $t = 0, 1, \dots, t_0 - 1$ **do**

$$\begin{aligned} \mathbf{U}^{t+1} &= \mathbf{U}^t - \eta_t \nabla_{\mathbf{U}} g(\mathbf{U}^t, \mathbf{V}^t, \mathbf{W}^t), \\ \mathbf{V}^{t+1} &= \mathbf{V}^t - \eta_t \nabla_{\mathbf{V}} g(\mathbf{U}^t, \mathbf{V}^t, \mathbf{W}^t), \\ \mathbf{W}^{t+1} &= \mathbf{W}^t - \eta_t \nabla_{\mathbf{W}} g(\mathbf{U}^t, \mathbf{V}^t, \mathbf{W}^t), \end{aligned}$$

where the gradients are given in (278).

Algorithm 8 Spectral initialization for nonconvex tensor completion (asymmetric case)

- 1: Let $\mathbf{U}\mathbf{A}\mathbf{U}^\top$ be the rank- r eigen-decomposition of $\mathcal{P}_{\text{off-diag}}(\mathbf{A}\mathbf{A}^\top)$ where $\mathbf{A} = \text{unfold}(\mathbf{T})$ is the mode-1 matricization of \mathbf{T} , and $\mathcal{P}_{\text{off-diag}}(\mathbf{Z})$ extracts out the off-diagonal entries of \mathbf{Z} .
 - 2: **Output:** initial estimates $\mathbf{U}^0 \in \mathbb{R}^{d_1 \times r}$, $\mathbf{V}^0 \in \mathbb{R}^{d_2 \times r}$, $\mathbf{W}^0 \in \mathbb{R}^{d_3 \times r}$ on the basis of $\mathbf{U} \in \mathbb{R}^{d_1 \times r}$ using Algorithm 9.
-

To begin with, let us generate the true tensor $\mathbf{T}^* = \sum_{1 \leq i \leq r} \mathbf{u}_i^* \otimes \mathbf{v}_i^* \otimes \mathbf{w}_i^*$ via the following procedure: (1) generate $\hat{\mathbf{u}}_i^* \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_1})$, $\hat{\mathbf{v}}_i^* \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_2})$ and $\hat{\mathbf{w}}_i^* \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_3})$, and (2) set $\lambda_i^* := \|\hat{\mathbf{u}}_i^*\|_2 \|\hat{\mathbf{v}}_i^*\|_2 \|\hat{\mathbf{w}}_i^*\|_2$, $\mathbf{u}_i^* = \lambda_i^{*1/3} \hat{\mathbf{u}}_i^*$, $\mathbf{v}_i^* = \lambda_i^{*1/3} \hat{\mathbf{v}}_i^*$ and $\mathbf{w}_i^* = \lambda_i^{*1/3} \hat{\mathbf{w}}_i^*$. Akin to the symmetric case, we choose the algorithmic parameters to be $L = r^2$, $\epsilon_{\text{th}} = 0.4$ and $t_0 = 100$. The noise components are generated as i.i.d. Gaussians, namely, $E_{i,j,k} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$. The signal-to-noise ratio (SNR) is defined as $\text{SNR} = \|\mathbf{T}^*\|_{\text{F}}^2 / (\sigma^2 d_1 d_2 d_3)$, and $\text{SNR} = \infty$ stands for the noiseless case (i.e. the case with $\sigma = 0$). The stepsize is set to be $\eta_t \equiv \eta = 1 / (2 \max_i \|\mathbf{u}_i^0\|_2^{4/3})$, where $\mathbf{U}^0 = [\mathbf{u}_1^0, \dots, \mathbf{u}_r^0]$ is the initial estimate generated by Algorithm 8. Figure 2 illustrates the numerical convergence rates of the proposed algorithm, where we set $d_1 = 100$, $d_2 = 150$, $d_3 = 200$, $r = 4$ and $p = 0.05$. Since one can only recover the tensor factors up to global signs and global permutation, the relative estimation errors of \mathbf{U}^t are defined to be

$$\begin{aligned} & \min_{\mathbf{\Pi} \in \text{perm}_r, \mathbf{S} = \text{diag}(s_i) \in \mathbb{R}^{r \times r}, s_i = \pm 1} \left\| \mathbf{U}^t \mathbf{\Pi} \mathbf{S} - \mathbf{U}^* \right\|_{\text{F}} / \|\mathbf{U}^*\|_{\text{F}}, \\ & \min_{\mathbf{\Pi} \in \text{perm}_r, \mathbf{S} = \text{diag}(s_i) \in \mathbb{R}^{r \times r}, s_i = \pm 1} \left\| \mathbf{U}^t \mathbf{\Pi} \mathbf{S} - \mathbf{U}^* \right\|_{2, \infty} / \|\mathbf{U}^*\|_{2, \infty}, \end{aligned}$$

where perm_r stands for the set of $r \times r$ permutation matrices. The error metrics for \mathbf{V}^t and \mathbf{W}^t can be defined analogously. The relative Euclidean and $\ell_{2, \infty}$ estimation errors of \mathbf{T} are defined as $\|\mathbf{T}^t - \mathbf{T}^*\|_{\text{F}} / \|\mathbf{T}^*\|_{\text{F}}$ and $\|\mathbf{T}^t - \mathbf{T}^*\|_{2, \infty} / \|\mathbf{T}^*\|_{2, \infty}$, respectively, where $\mathbf{T}^t = \sum_{1 \leq i \leq r} \mathbf{u}_i^t \otimes \mathbf{v}_i^t \otimes \mathbf{w}_i^t$. As can be seen from the plots, the estimation errors decay geometrically fast in the noiseless case. In the noisy case, the numerical estimation errors of the algorithm also converge geometrically fast until an error floor is hit.

Moving beyond the above synthetic data, we apply our methods to an simulated MRI brain image dataset [CKK⁺97], which is available online at the McGill University McConnell Brain Imaging Centre and has also been studied in prior work [XYZ17]. The database consists of pre-computed simulated brain data, and a set of parameters can be set to generate the data accordingly. In this series of experiments, we choose the parameters to be T1 modality, 1mm slice thickness, 1% noise and 20% RF. The resulting data is a three-order tensor in $\mathbb{R}^{181 \times 217 \times 181}$, where each slice in any mode corresponds to a brain image. We use the tensor decomposition algorithm in [AGJ15] to decompose the original data, and keep the top-36 components whose energy accounts for 85% of the original tensor. Hence, the resulting low-rank tensor \mathbf{T}^* preserves the major information of the original tensor, and is used as the ground truth in the simulation. We sample each entry of \mathbf{T}^* independently with probability p , and then inject i.i.d. Gaussian noise $\mathcal{N}(0, \sigma^2)$ to each observed entry (so as to simulate the noisy scenario). Based on the incomplete set of noisy samples, we apply the proposed algorithm to reconstruct \mathbf{T}^* ; the algorithmic parameters are chosen as above, and we denote by \mathbf{T} the resulting tensor estimate. Table 2 reports the relative errors $\|\mathbf{T} - \mathbf{T}^*\|_{\text{F}} / \|\mathbf{T}^*\|_{\text{F}}$ and $\|\mathbf{T} - \mathbf{T}^*\|_{\infty} / \|\mathbf{T}^*\|_{\infty}$. As can be seen in the table, the performance of our algorithm is quite favorable for low-rank tensor reconstruction. In particular, given that these real data tensor examples might not satisfy

Algorithm 9 Retrieval of low-rank tensor factors from a given subspace estimate (asymmetric case)

- 1: **Input:** number of restarts L , pruning threshold ϵ_{th} , subspace estimate $\mathbf{U} \in \mathbb{R}^{d_1 \times r}$ given by Algorithm 8.
- 2: **for** $\tau = 1, \dots, L$ **do**
- 3: Generate an independent Gaussian vector $\mathbf{g}^\tau \sim \mathcal{N}(0, \mathbf{I}_{d_1})$.
- 4:

$$(\boldsymbol{\nu}^{(1),\tau}, \boldsymbol{\nu}^{(2),\tau}, \boldsymbol{\nu}^{(3),\tau}, \lambda_\tau, \text{gap}_\tau) \leftarrow \text{RETRIEVE-ONE-TENSOR-FACTOR-ASYM}(\mathbf{T}, p, \mathbf{U}, \mathbf{g}^\tau).$$

- 5: Generate tensor factor estimates

$$\{(\mathbf{u}^1, \mathbf{v}^1, \mathbf{w}^1, \lambda_1), \dots, (\mathbf{u}^\tau, \mathbf{v}^\tau, \mathbf{w}^\tau, \lambda_\tau)\} \leftarrow \text{PRUNE-ASYM}(\{(\boldsymbol{\nu}^{(1),\tau}, \boldsymbol{\nu}^{(2),\tau}, \boldsymbol{\nu}^{(3),\tau}, \lambda_\tau, \text{gap}_\tau)\}_{\tau=1}^L, \epsilon_{\text{th}}).$$

- 6: **Output:** initial estimate

$$\mathbf{U}^0 = [\lambda_1^{1/3} \mathbf{u}^1, \dots, \lambda_r^{1/3} \mathbf{u}^r], \quad \mathbf{V}^0 = [\lambda_1^{1/3} \mathbf{v}^1, \dots, \lambda_r^{1/3} \mathbf{v}^r], \quad \text{and} \quad \mathbf{W}^0 = [\lambda_1^{1/3} \mathbf{w}^1, \dots, \lambda_r^{1/3} \mathbf{w}^r].$$

- 1: **function** RETRIEVE-ONE-TENSOR-FACTOR-ASYM($\mathbf{T}, p, \mathbf{U}, \mathbf{g}$)
- 2: Compute

$$\boldsymbol{\theta} = \mathbf{U}\mathbf{U}^\top \mathbf{g} =: \mathcal{P}_{\mathbf{U}}(\mathbf{g}), \tag{279a}$$

$$\mathbf{M} = p^{-1} \mathbf{T} \times_1 \boldsymbol{\theta}, \tag{279b}$$

where \times_1 is defined in Section 2.4.

- 3: Let $\boldsymbol{\nu}^{(2)}$ (resp. $\boldsymbol{\nu}^{(3)}$) be the leading left (resp. right) singular vector of \mathbf{M} . Let $\boldsymbol{\nu}^{(1)} = p^{-1} \mathbf{T} \times_2 \boldsymbol{\nu}^{(2)} \times_3 \boldsymbol{\nu}^{(3)}$ and set $\lambda = \|\boldsymbol{\nu}^{(1)}\|_2^{1/3}$.
 - 4: **return** $(\boldsymbol{\nu}^{(1)}, \boldsymbol{\nu}^{(2)}, \boldsymbol{\nu}^{(3)}, \lambda, \sigma_1(\mathbf{M}) - \sigma_2(\mathbf{M}))$.
-

the assumptions imposed in our main theorems, these numerical experiments hint at the applicability of our algorithm to a broader set of problems.

E.4 Analysis ideas

Before describing the proof ideas, we define the following incoherence parameters and condition number that, similar to the symmetric case, play a crucial role in our theoretical development.

Definition E.1. Define the incoherence parameters and the condition number of \mathbf{T}^* as follows

$$\mu_0 := \frac{d_1 d_2 d_3 \|\mathbf{T}^*\|_\infty^2}{\|\mathbf{T}^*\|_{\mathbb{F}}^2}, \tag{280a}$$

$$\mu_1 := \max \left\{ \frac{d_1 \|\mathbf{u}_i^*\|_\infty^2}{\|\mathbf{u}_i^*\|_2^2}, \frac{d_2 \|\mathbf{v}_i^*\|_\infty^2}{\|\mathbf{v}_i^*\|_2^2}, \frac{d_3 \|\mathbf{w}_i^*\|_\infty^2}{\|\mathbf{w}_i^*\|_2^2} \right\}, \tag{280b}$$

$$\mu_2 := \max \left\{ \frac{d_1 \langle \mathbf{u}_i^*, \mathbf{u}_j^* \rangle^2}{\|\mathbf{u}_i^*\|_2^2 \|\mathbf{u}_j^*\|_2^2}, \frac{d_2 \langle \mathbf{v}_i^*, \mathbf{v}_j^* \rangle^2}{\|\mathbf{v}_i^*\|_2^2 \|\mathbf{v}_j^*\|_2^2}, \frac{d_3 \langle \mathbf{w}_i^*, \mathbf{w}_j^* \rangle^2}{\|\mathbf{w}_i^*\|_2^2 \|\mathbf{w}_j^*\|_2^2} \right\}, \tag{280c}$$

$$\kappa := \frac{\max_i \{\|\mathbf{u}_i^*\|_2 \|\mathbf{v}_i^*\|_2 \|\mathbf{w}_i^*\|_2\}}{\min_i \{\|\mathbf{u}_i^*\|_2 \|\mathbf{v}_i^*\|_2 \|\mathbf{w}_i^*\|_2\}}. \tag{280d}$$

For notational convenience, we shall set

$$\mu := \max \{\mu_0, \mu_1, \mu_2\}, \quad d_{\min} := \min \{d_1, d_2, d_3\} \quad \text{and} \quad d_{\max} := \max \{d_1, d_2, d_3\}. \tag{280e}$$

As has been made clear in the symmetric case, at the heart of our analysis lie two crucial components: (1) the geometric property of a noiseless version of the loss function, and (2) reasonably well initial estimates

```

1: function PRUNE-ASYM( $\{(\boldsymbol{\nu}^{(1),\tau}, \boldsymbol{\nu}^{(2),\tau}, \boldsymbol{\nu}^{(3),\tau}, \lambda_\tau, \text{gap}_\tau)\}_{\tau=1}^L, \epsilon_{\text{th}})$ )
2:   Set  $\Theta = \{(\boldsymbol{\nu}^{(1),\tau}, \boldsymbol{\nu}^{(2),\tau}, \boldsymbol{\nu}^{(3),\tau}, \lambda_\tau, \text{gap}_\tau)\}_{\tau=1}^L$ .
3:   for  $i = 1, \dots, r$  do
4:     Choose  $(\boldsymbol{\nu}^{(1),\tau}, \boldsymbol{\nu}^{(2),\tau}, \boldsymbol{\nu}^{(3),\tau}, \lambda_\tau, \text{gap}_\tau)$  from  $\Theta$  with the largest  $\text{gap}_\tau$ ; set  $\mathbf{u}^i = \boldsymbol{\nu}^{(1),\tau}, \mathbf{v}^i = \boldsymbol{\nu}^{(2),\tau}, \mathbf{w}^i = \boldsymbol{\nu}^{(3),\tau}$  and  $\lambda_i = \lambda_\tau$ .
5:     Update  $\Theta \leftarrow \Theta \setminus \{(\boldsymbol{\nu}^{(1),\tau}, \boldsymbol{\nu}^{(2),\tau}, \boldsymbol{\nu}^{(3),\tau}, \lambda_\tau, \text{gap}_\tau) \in \Theta : |\langle \boldsymbol{\nu}^{(2),\tau}, \mathbf{v}^i \rangle| > 1 - \epsilon_{\text{th}}\}$ .
6:   return  $\{(\mathbf{u}^1, \mathbf{v}^1, \mathbf{w}^1, \lambda_1), \dots, (\mathbf{u}^r, \mathbf{v}^r, \mathbf{w}^r, \lambda_r)\}$ .

```

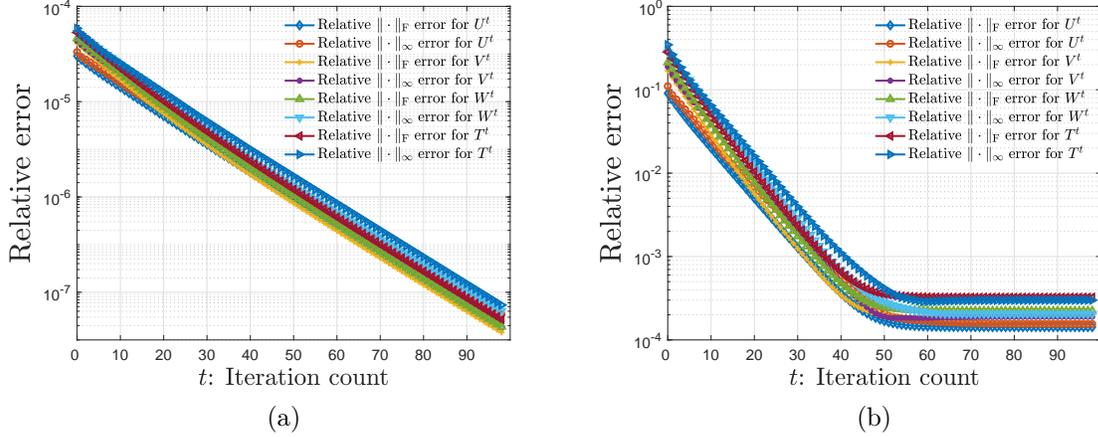


Figure 2: (a) relative errors of the estimates $\mathbf{U}^t, \mathbf{V}^t, \mathbf{W}^t$ and \mathbf{T}^t vs. iteration count t for noiseless tensor completion, where $d_1 = 100, d_2 = 150, d_3 = 200, r = 4, p = 0.05$; (b) relative errors of the estimates $\mathbf{U}^t, \mathbf{V}^t, \mathbf{W}^t$ and \mathbf{T}^t vs. iteration count t for noisy tensor completion, where $d_1 = 100, d_2 = 150, d_3 = 200, r = 4, p = 0.05, \text{SNR} = 10$.

for tensor factors (in the entrywise sense). Rather than providing a complete analysis for the asymmetric case (which would be very long), we shall only point out the important steps needed to extend these two parts for the asymmetric case.

1. Local optimization landscape. Similar to the symmetric counterpart, the key step of the local convergence analysis lies in establishing the favorable geometric property (i.e. local strong convexity and smoothness) of the following noiseless regularized loss function

$$g_{\text{clean}}(\mathbf{U}, \mathbf{V}, \mathbf{W}) = \frac{1}{6p} \left\| \mathcal{P}_\Omega \left(\sum_{i=1}^r \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i - \mathbf{T}^* \right) \right\|_{\text{F}}^2 + \text{reg}(\mathbf{U}, \mathbf{V}, \mathbf{W}), \quad (281)$$

where the regularization term $\text{reg}(\mathbf{U}, \mathbf{V}, \mathbf{W})$ is defined in (277). In words, this is a simplified version of the original loss function (276) by dropping the influence of the noise. Lemma E.2 below demonstrates that g_{clean} is locally strongly convex and smooth in the neighborhood of the ground truth.

Lemma E.2 (Local strong convexity and smoothness). *Suppose that*

$$p \geq c_0 \max \left\{ \frac{\log^3 d_{\max}}{\sqrt{d_1 d_2 d_3}}, \frac{\mu^2 r^2 d_{\max} \log^5 d_{\max}}{d_1 d_2 d_3} \right\}, \quad r \leq c_1 \sqrt{\frac{d_{\min}}{\mu}} \quad (282)$$

and that the regularization parameter obeys $|\alpha_i - \lambda_i^{*2/3}| \leq c_2 \lambda_{\min}^{*2/3}$ for some sufficiently large (resp. small)

Table 2: Relative ℓ_2 and $\ell_{2,\infty}$ errors for varying p and SNR in the MRI data experiments.

(p, SNR)	$\ \mathbf{T} - \mathbf{T}^*\ _{\text{F}}/\ \mathbf{T}^*\ _{\text{F}}$	$\ \mathbf{T} - \mathbf{T}^*\ _{\infty}/\ \mathbf{T}^*\ _{\infty}$
(0.05, 10)	0.2231	0.8774
(0.05, ∞)	0.2149	0.7856
(0.1, 10)	0.1472	0.5343
(0.1, ∞)	0.1398	0.3991
(0.2, 10)	0.0841	0.1417
(0.2, ∞)	0.0686	0.1193

constant $c_0 > 0$ (resp. $c_1, c_2 > 0$). Then with probability greater than $1 - O(d_{\min}^{-10})$,

$$\text{vec} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \\ \mathbf{Z} \end{pmatrix}^{\top} \nabla^2 g_{\text{clean}}(\mathbf{U}, \mathbf{V}, \mathbf{W}) \text{vec} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \\ \mathbf{Z} \end{pmatrix} \leq 4\lambda_{\max}^{*4/3} \left(\|\mathbf{X}\|_{\text{F}}^2 + \|\mathbf{Y}\|_{\text{F}}^2 + \|\mathbf{Z}\|_{\text{F}}^2 \right), \quad (283)$$

$$\text{vec} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \\ \mathbf{Z} \end{pmatrix}^{\top} \nabla^2 g_{\text{clean}}(\mathbf{U}, \mathbf{V}, \mathbf{W}) \text{vec} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \\ \mathbf{Z} \end{pmatrix} \geq \frac{1}{2}\lambda_{\min}^{*4/3} \left(\|\mathbf{X}\|_{\text{F}}^2 + \|\mathbf{Y}\|_{\text{F}}^2 + \|\mathbf{Z}\|_{\text{F}}^2 \right) \quad (284)$$

holds simultaneously for all $\mathbf{X} \in \mathbb{R}^{d_1 \times r}$, $\mathbf{Y} \in \mathbb{R}^{d_2 \times r}$, $\mathbf{Z} \in \mathbb{R}^{d_3 \times r}$ and all $\mathbf{U} \in \mathbb{R}^{d_1 \times r}$, $\mathbf{V} \in \mathbb{R}^{d_2 \times r}$, $\mathbf{W} \in \mathbb{R}^{d_3 \times r}$ obeying

$$\|\mathbf{U} - \mathbf{U}^* \mathbf{S}^{(1)}\|_{\text{F}} \leq \delta \|\mathbf{U}^*\|_{\text{F}} \quad \text{and} \quad \left\| \mathbf{U} - \mathbf{U}^* \mathbf{S}^{(1)} \right\|_{2,\infty} \leq \delta \|\mathbf{U}^*\|_{2,\infty}; \quad (285\text{a})$$

$$\|\mathbf{V} - \mathbf{V}^* \mathbf{S}^{(2)}\|_{\text{F}} \leq \delta \|\mathbf{V}^*\|_{\text{F}} \quad \text{and} \quad \left\| \mathbf{V} - \mathbf{V}^* \mathbf{S}^{(2)} \right\|_{2,\infty} \leq \delta \|\mathbf{V}^*\|_{2,\infty}; \quad (285\text{b})$$

$$\|\mathbf{W} - \mathbf{W}^* \mathbf{S}^{(3)}\|_{\text{F}} \leq \delta \|\mathbf{W}^*\|_{\text{F}} \quad \text{and} \quad \left\| \mathbf{W} - \mathbf{W}^* \mathbf{S}^{(3)} \right\|_{2,\infty} \leq \delta \|\mathbf{W}^*\|_{2,\infty}. \quad (285\text{c})$$

Here, $\delta \leq c_3/(\mu^{3/2}r)$ for some sufficiently small constant $c_3 > 0$, and $\mathbf{S}^{(1)}$, $\mathbf{S}^{(2)}$ and $\mathbf{S}^{(3)}$ are some diagonal matrices in $\mathbb{R}^{r \times r}$ such that for each $1 \leq i \leq r$, two of $S_{i,i}^{(1)}$, $S_{i,i}^{(2)}$ and $S_{i,i}^{(3)}$ equal to -1 with the remaining one equal to 1.

Proof. See Appendix E.5. ■

In a nutshell, Lemma E.2 confirms local strong convexity and smoothness of the noiseless regularized loss $g_{\text{clean}}(\mathbf{U}, \mathbf{V}, \mathbf{W})$, provided that (282) holds and that the matrices \mathbf{U} , \mathbf{V} and \mathbf{W} are sufficiently close to the ground truth in every single row. This is similar to the property of the symmetric counterpart $f_{\text{clean}}(\mathbf{U})$ (cf. Lemma 5.1 in Appendix 5.1), except that we need to deal with asymmetric tensor factors as well as additional regularization terms. In particular, when $d_1 \asymp d_2 \asymp d_3 \asymp d$ and $\mu, r \asymp 1$, Condition (282) reduces to $p \gg d^{-3/2} \log^3 d$ and $r \ll \sqrt{d}$, which resembles (42) in Lemma 5.1 derived for the symmetric case.

Having established the preceding local geometric properties of g_{clean} , one can then argue similarly as in Lemmas 5.3 and 5.6 in Appendix 5.1 to prove that gradient descent converges linearly, as long as it is provided with an initial estimate satisfying the condition (285). Here, we emphasize that the regularization term (277), which essentially balances the sizes of the three tensor factors, is crucial for the local strong convexity and smoothness of g_{clean} to hold.

2. Guaranteeing a reasonably good initialization. Another crucial ingredient lies in guaranteeing an initial estimate with sufficiently good accuracy. Recall that our initialization scheme consists of two stages: (1) subspace estimation, and (2) retrieval of individual tensor factors.

- The subspace estimation part remains largely unchanged: we shall unfold the observed tensor along the 1-st mode to estimate the subspace spanned by tensor factors $\{\mathbf{u}_i^*\}_{i=1}^r$, and the $\ell_{2,\infty}$ subspace estimation accuracy can be established by invoking the main theorems of our companion paper [CLC+20].

- Regarding the retrieval of individual tensor factors, the key observation is that: the random vector \mathbf{g}^τ we generate satisfies

$$\mathbf{T}^* \times_1 \mathcal{P}_{\mathbf{U}^*}(\mathbf{g}^\tau) = \sum_{i=1}^r \lambda_i^* \langle \bar{\mathbf{u}}_i^*, \mathcal{P}_{\mathbf{U}^*}(\mathbf{g}^\tau) \rangle \bar{\mathbf{v}}_i^* \bar{\mathbf{w}}_i^{*\top} = \sum_{i=1}^r \lambda_i^* \langle \bar{\mathbf{u}}_i^*, \mathbf{g}^\tau \rangle \bar{\mathbf{v}}_i^* \bar{\mathbf{w}}_i^{*\top},$$

where $\mathcal{P}_{\mathbf{U}^*}$ is the projection onto the subspace spanned by $\{\mathbf{u}_i^*\}_{i=1}^r$, and

$$\bar{\mathbf{u}}_i^* = \frac{1}{\|\mathbf{u}_i^*\|_2} \mathbf{u}_i^*, \quad \bar{\mathbf{v}}_i^* = \frac{1}{\|\mathbf{v}_i^*\|_2} \mathbf{v}_i^*, \quad \text{and} \quad \bar{\mathbf{w}}_i^* = \frac{1}{\|\mathbf{w}_i^*\|_2} \mathbf{w}_i^*. \quad (286)$$

Given a sufficiently accurate subspace estimate \mathbf{U} for \mathbf{U}^* obtained in the subspace estimation stage, for each $1 \leq i \leq r$, there exists at least a point $1 \leq \tau \leq L$ such that the spectral gap of the population version of $\mathbf{T} \times_1 \boldsymbol{\theta}^\tau = \mathbf{T} \times_1 \mathcal{P}_{\mathbf{U}}(\mathbf{g}^\tau)$ (with respect to the sampling and noise) is large enough and that the perturbation is sufficiently small. As a result, the top left (resp. right) singular vector $\boldsymbol{\nu}^{(2)}$ (resp. $\boldsymbol{\nu}^{(3)}$) of $\mathbf{T} \times_1 \boldsymbol{\theta}^\tau$ is close to the (normalized) tensor factor $\bar{\mathbf{v}}_i^*$ (resp. $\bar{\mathbf{w}}_i^*$) both in the ℓ_2 and ℓ_∞ norm (similar to Lemmas 5.13-5.14 in Section 5.3 for the symmetric case). In turn, this further allows us to reliably estimate $\bar{\mathbf{u}}_i^*$ and the magnitude λ_i^* , in a way similar to what we have done in Lemma 5.16 in Section 5.3. By repeating the procedures with random restarts and invoking a pruning procedure similar to the symmetric case, we can hope to recover all tensor factors with high probability.

Following the above strategies, one could adapt the proofs of Theorems 5.9-5.10 in Section 5.3 to show that: our initial estimates $\{\mathbf{u}_i, \mathbf{v}_i, \mathbf{w}_i\}_{i=1}^r$ are all exceedingly close to the ground truth in the entrywise sense (up to global permutation and global signs). This in turn confirms that the algorithm will enter a locally strongly convex and smooth region as characterized in Lemma E.2, thus leading to our performance guarantees for the entire algorithm. Once again, while the analysis ideas for the asymmetric case bear much resemblance to the symmetric counterpart, a complete proof has to be fairly long due to more clumsy notation compared to the symmetric case; for the sake of brevity, we do not provide the full proof here.

E.5 Proof of Lemma E.2

Fix arbitrary matrices $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_r] \in \mathbb{R}^{d_1 \times r}$, $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_r] \in \mathbb{R}^{d_2 \times r}$ and $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_r] \in \mathbb{R}^{d_3 \times r}$. Direct computation reveals that

$$\begin{aligned} & \text{vec} \left(\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \\ \mathbf{Z} \end{bmatrix} \right)^\top \nabla^2 g_{\text{clean}}(\mathbf{U}, \mathbf{V}, \mathbf{W}) \text{vec} \left(\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \\ \mathbf{Z} \end{bmatrix} \right) \\ &= \frac{1}{3p} \left\| \mathcal{P}_\Omega \left(\sum_{i=1}^r \mathbf{x}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i + \mathbf{u}_i \otimes \mathbf{y}_i \otimes \mathbf{w}_i + \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{z}_i \right) \right\|_{\text{F}}^2 \\ &+ \frac{2}{3p} \left\langle \sum_{i=1}^r \mathbf{x}_i \otimes \mathbf{y}_i \otimes \mathbf{w}_i + \mathbf{x}_i \otimes \mathbf{v}_i \otimes \mathbf{z}_i + \mathbf{u}_i \otimes \mathbf{y}_i \otimes \mathbf{z}_i, \mathcal{P}_\Omega \left(\sum_{i=1}^r \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i - \mathbf{T}^* \right) \right\rangle \\ &+ \frac{1}{3} \sum_{i=1}^r \alpha_i (\langle \mathbf{x}_i, \mathbf{u}_i \rangle - \langle \mathbf{y}_i, \mathbf{v}_i \rangle)^2 + \frac{1}{6} \sum_{i=1}^r \alpha_i (\|\mathbf{u}_i\|_2^2 - \|\mathbf{v}_i\|_2^2) (\|\mathbf{x}_i\|_2^2 - \|\mathbf{y}_i\|_2^2) \\ &+ \frac{1}{3} \sum_{i=1}^r \alpha_i (\langle \mathbf{x}_i, \mathbf{u}_i \rangle - \langle \mathbf{z}_i, \mathbf{w}_i \rangle)^2 + \frac{1}{6} \sum_{i=1}^r \alpha_i (\|\mathbf{u}_i\|_2^2 - \|\mathbf{w}_i\|_2^2) (\|\mathbf{x}_i\|_2^2 - \|\mathbf{z}_i\|_2^2) \\ &+ \frac{1}{3} \sum_{i=1}^r \alpha_i (\langle \mathbf{y}_i, \mathbf{v}_i \rangle - \langle \mathbf{z}_i, \mathbf{w}_i \rangle)^2 + \frac{1}{6} \sum_{i=1}^r \alpha_i (\|\mathbf{v}_i\|_2^2 - \|\mathbf{w}_i\|_2^2) (\|\mathbf{y}_i\|_2^2 - \|\mathbf{z}_i\|_2^2). \end{aligned}$$

For notational convenience, we shall define

$$\boldsymbol{\Delta}^{\mathbf{U}} = [\boldsymbol{\Delta}_1^{\mathbf{U}}, \dots, \boldsymbol{\Delta}_r^{\mathbf{U}}] := \mathbf{U} \mathbf{S}^{(1)} - \mathbf{U}^* \in \mathbb{R}^{d_1 \times r},$$

and define $\Delta^V \in \mathbb{R}^{d_2 \times r}$ and $\Delta^W \in \mathbb{R}^{d_3 \times r}$ in an analogous manner. Without loss of generality, it is assumed that $\mathbf{S}^{(1)} = \mathbf{I}_r$ and $\mathbf{S}^{(2)} = \mathbf{S}^{(3)} = -\mathbf{I}_r$.

With the above notation in place, one can decompose

$$\text{vec} \left(\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \\ \mathbf{Z} \end{bmatrix} \right)^\top \nabla^2 g_{\text{clean}}(\mathbf{U}, \mathbf{V}, \mathbf{W}) \text{vec} \left(\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \\ \mathbf{Z} \end{bmatrix} \right) = \sum_{i=1}^4 \beta_i$$

where $\beta_i, 1 \leq i \leq 4$ are given respectively by

$$\begin{aligned} \beta_1 &:= \frac{1}{3} \left\| \sum_{i=1}^r \mathbf{x}_i \otimes (-\mathbf{v}_i^*) \otimes (-\mathbf{w}_i^*) + \mathbf{u}_i^* \otimes \mathbf{y}_i \otimes (-\mathbf{w}_i^*) + \mathbf{u}_i^* \otimes (-\mathbf{v}_i^*) \otimes \mathbf{z}_i \right\|_{\mathbb{F}}^2 \\ &\quad + \frac{1}{3} \sum_{i=1}^r \lambda_i^{*2/3} \left[(\langle \mathbf{x}_i, \mathbf{u}_i^* \rangle - \langle \mathbf{y}_i, -\mathbf{v}_i^* \rangle)^2 + (\langle \mathbf{x}_i, \mathbf{u}_i^* \rangle - \langle \mathbf{z}_i, -\mathbf{w}_i^* \rangle)^2 + (\langle \mathbf{y}_i, -\mathbf{v}_i^* \rangle - \langle \mathbf{z}_i, -\mathbf{w}_i^* \rangle)^2 \right], \end{aligned} \quad (287a)$$

$$\begin{aligned} \beta_2 &:= \frac{1}{3p} \left\| \mathcal{P}_\Omega \left(\sum_{i=1}^r \mathbf{x}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i + \mathbf{u}_i \otimes \mathbf{y}_i \otimes \mathbf{w}_i + \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{z}_i \right) \right\|_{\mathbb{F}}^2 \\ &\quad - \frac{1}{3} \left\| \sum_{i=1}^r \mathbf{x}_i \otimes (-\mathbf{v}_i^*) \otimes (-\mathbf{w}_i^*) + \mathbf{u}_i^* \otimes \mathbf{y}_i \otimes (-\mathbf{w}_i^*) + \mathbf{u}_i^* \otimes (-\mathbf{v}_i^*) \otimes \mathbf{z}_i \right\|_{\mathbb{F}}^2, \end{aligned} \quad (287b)$$

$$\beta_3 := \frac{2}{3p} \left\langle \sum_{i=1}^r \mathbf{x}_i \otimes \mathbf{y}_i \otimes \mathbf{w}_i + \mathbf{x}_i \otimes \mathbf{v}_i \otimes \mathbf{z}_i + \mathbf{u}_i \otimes \mathbf{y}_i \otimes \mathbf{z}_i, \mathcal{P}_\Omega \left(\sum_{i=1}^r \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i - \mathbf{T}^* \right) \right\rangle, \quad (287c)$$

$$\begin{aligned} \beta_4 &:= \frac{1}{3} \sum_{i=1}^r (\alpha_i - \lambda_i^{*2/3}) \left[(\langle \mathbf{x}_i, \mathbf{u}_i^* \rangle - \langle \mathbf{y}_i, -\mathbf{v}_i^* \rangle)^2 + (\langle \mathbf{x}_i, \mathbf{u}_i^* \rangle - \langle \mathbf{z}_i, -\mathbf{w}_i^* \rangle)^2 + (\langle \mathbf{y}_i, -\mathbf{v}_i^* \rangle - \langle \mathbf{z}_i, -\mathbf{w}_i^* \rangle)^2 \right] \\ &\quad + \frac{1}{3} \sum_{i=1}^r \alpha_i \left[(\langle \mathbf{x}_i, \Delta_i^U \rangle - \langle \mathbf{y}_i, \Delta_i^V \rangle)^2 + (\langle \mathbf{x}_i, \Delta_i^U \rangle - \langle \mathbf{z}_i, \Delta_i^W \rangle)^2 + (\langle \mathbf{y}_i, \Delta_i^V \rangle - \langle \mathbf{z}_i, \Delta_i^W \rangle)^2 \right] \\ &\quad + \frac{2}{3} \sum_{i=1}^r \alpha_i \left[(\langle \mathbf{x}_i, \mathbf{u}_i^* \rangle - \langle \mathbf{y}_i, -\mathbf{v}_i^* \rangle) (\langle \mathbf{x}_i, \Delta_i^U \rangle - \langle \mathbf{y}_i, \Delta_i^V \rangle) + (\langle \mathbf{x}_i, \mathbf{u}_i^* \rangle - \langle \mathbf{z}_i, -\mathbf{w}_i^* \rangle) (\langle \mathbf{x}_i, \Delta_i^U \rangle - \langle \mathbf{z}_i, \Delta_i^W \rangle) \right. \\ &\quad \left. + (\langle \mathbf{y}_i, -\mathbf{v}_i^* \rangle - \langle \mathbf{z}_i, -\mathbf{w}_i^* \rangle) (\langle \mathbf{y}_i, \Delta_i^V \rangle - \langle \mathbf{z}_i, \Delta_i^W \rangle) \right] \\ &\quad + \frac{1}{6} \sum_{i=1}^r \alpha_i \left[(\|\mathbf{u}_i\|_2^2 - \|\mathbf{v}_i\|_2^2) (\|\mathbf{x}_i\|_2^2 - \|\mathbf{y}_i\|_2^2) + (\|\mathbf{u}_i\|_2^2 - \|\mathbf{w}_i\|_2^2) (\|\mathbf{x}_i\|_2^2 - \|\mathbf{z}_i\|_2^2) \right. \\ &\quad \left. + (\|\mathbf{v}_i\|_2^2 - \|\mathbf{w}_i\|_2^2) (\|\mathbf{y}_i\|_2^2 - \|\mathbf{z}_i\|_2^2) \right]. \end{aligned} \quad (287d)$$

In what follows, we shall demonstrate that β_1 is the dominant term, with the remaining terms being negligible compared to β_1 . Here, we note that the proof idea is almost identical to that of the symmetric case (cf. Lemma 5.1 in Appendix 5.1). For the sake of conciseness, we will focus only on the part where the symmetric and the asymmetric cases differ, and omit the proof details when their analyses are similar.

Bounding β_1 Let us first expand

$$\begin{aligned} \beta_1 &= \frac{1}{3} \left\| \sum_{i=1}^r \mathbf{x}_i \otimes \mathbf{v}_i^* \otimes \mathbf{w}_i^* - \mathbf{u}_i^* \otimes \mathbf{y}_i \otimes \mathbf{w}_i^* - \mathbf{u}_i^* \otimes \mathbf{v}_i^* \otimes \mathbf{z}_i \right\|_{\mathbb{F}}^2 \\ &\quad + \frac{1}{3} \sum_{i=1}^r \lambda_i^{*2/3} \left[(\langle \mathbf{x}_i, \mathbf{u}_i^* \rangle + \langle \mathbf{y}_i, \mathbf{v}_i^* \rangle)^2 + (\langle \mathbf{x}_i, \mathbf{u}_i^* \rangle + \langle \mathbf{z}_i, \mathbf{w}_i^* \rangle)^2 + (\langle \mathbf{y}_i, \mathbf{v}_i^* \rangle - \langle \mathbf{z}_i, \mathbf{w}_i^* \rangle)^2 \right] \\ &= \underbrace{\frac{1}{3} \left\| \sum_{i=1}^r \mathbf{x}_i \otimes \mathbf{v}_i^* \otimes \mathbf{w}_i^* \right\|_{\mathbb{F}}^2 + \frac{1}{3} \left\| \sum_{i=1}^r \mathbf{u}_i^* \otimes \mathbf{y}_i \otimes \mathbf{w}_i^* \right\|_{\mathbb{F}}^2 + \frac{1}{3} \left\| \sum_{i=1}^r \mathbf{u}_i^* \otimes \mathbf{v}_i^* \otimes \mathbf{z}_i \right\|_{\mathbb{F}}^2}_{=: \gamma_1} \end{aligned}$$

$$\begin{aligned}
& + \underbrace{\frac{2}{3} \sum_{i=1}^r \lambda_i^{*2/3} \langle \mathbf{x}_i, \mathbf{u}_i^* \rangle^2 + \frac{2}{3} \sum_{i=1}^r \lambda_i^{*2/3} \langle \mathbf{y}_i, \mathbf{v}_i^* \rangle^2 + \frac{2}{3} \sum_{i=1}^r \lambda_i^{*2/3} \langle \mathbf{z}_i, \mathbf{w}_i^* \rangle^2}_{=:\gamma_2} \\
& + \underbrace{-\frac{2}{3} \left\langle \sum_{i=1}^r \mathbf{x}_i \otimes \mathbf{v}_i^* \otimes \mathbf{w}_i^*, \sum_{i=1}^r \mathbf{u}_i^* \otimes \mathbf{y}_i \otimes \mathbf{w}_i^* \right\rangle + \frac{2}{3} \sum_{i=1}^r \lambda_i^{*2/3} \langle \mathbf{x}_i, \mathbf{u}_i^* \rangle \langle \mathbf{y}_i, \mathbf{v}_i^* \rangle}_{=:\gamma_3} \\
& + \underbrace{-\frac{2}{3} \left\langle \sum_{i=1}^r \mathbf{x}_i \otimes \mathbf{v}_i^* \otimes \mathbf{w}_i^*, \sum_{i=1}^r \mathbf{u}_i^* \otimes \mathbf{v}_i^* \otimes \mathbf{z}_i \right\rangle + \frac{2}{3} \sum_{i=1}^r \lambda_i^{*2/3} \langle \mathbf{x}_i, \mathbf{u}_i^* \rangle \langle \mathbf{z}_i, \mathbf{w}_i^* \rangle}_{=:\gamma_4} \\
& + \underbrace{\frac{2}{3} \left\langle \sum_{i=1}^r \mathbf{u}_i^* \otimes \mathbf{y}_i \otimes \mathbf{w}_i^*, \sum_{i=1}^r \mathbf{u}_i^* \otimes \mathbf{v}_i^* \otimes \mathbf{z}_i \right\rangle - \frac{2}{3} \sum_{i=1}^r \lambda_i^{*2/3} \langle \mathbf{y}_i, \mathbf{v}_i^* \rangle \langle \mathbf{z}_i, \mathbf{w}_i^* \rangle}_{=:\gamma_5},
\end{aligned}$$

leaving us with five terms to control.

1. Let us begin with γ_1 . Observe that we can express

$$\left\| \sum_{i=1}^r \mathbf{x}_i \otimes \mathbf{v}_i^* \otimes \mathbf{w}_i^* \right\|_{\mathbb{F}} = \|\mathbf{X}\mathbf{H}^\top\|_{\mathbb{F}}, \quad \text{with } \mathbf{H} = [\mathbf{v}_1^* \otimes \mathbf{w}_1^*, \dots, \mathbf{v}_r^* \otimes \mathbf{w}_r^*] \in \mathbb{R}^{d^2 \times r}$$

Arguing similarly as in the proof of Lemma D.1, one can derive

$$|\sigma_{\min}(\mathbf{H}) - \lambda_{\min}^{*2/3}| \ll \lambda_{\min}^{*2/3} \quad \text{and} \quad |\sigma_{\max}(\mathbf{H}) - \lambda_{\max}^{*2/3}| \ll \lambda_{\min}^{*2/3},$$

provided that $r \ll d_{\min}/\mu$. This leads to the following inequalities

$$\frac{19}{20} \lambda_{\min}^{*2/3} \|\mathbf{X}\|_{\mathbb{F}} \leq \sigma_{\min}(\mathbf{H}) \|\mathbf{X}\|_{\mathbb{F}} \leq \left\| \sum_{i=1}^r \mathbf{x}_i \otimes \mathbf{v}_i^* \otimes \mathbf{w}_i^* \right\|_{\mathbb{F}} \leq \sigma_{\max}(\mathbf{H}) \|\mathbf{X}\|_{\mathbb{F}} \leq \frac{11}{10} \lambda_{\max}^{*2/3} \|\mathbf{X}\|_{\mathbb{F}}.$$

Clearly, the same argument reveals that

$$\begin{aligned}
\frac{19}{20} \lambda_{\min}^{*2/3} \|\mathbf{Y}\|_{\mathbb{F}} & \leq \left\| \sum_{i=1}^r \mathbf{u}_i^* \otimes \mathbf{y}_i \otimes \mathbf{w}_i^* \right\|_{\mathbb{F}} \leq \frac{11}{10} \lambda_{\max}^{*2/3} \|\mathbf{Y}\|_{\mathbb{F}}, \\
\frac{19}{20} \lambda_{\min}^{*2/3} \|\mathbf{Z}\|_{\mathbb{F}} & \leq \left\| \sum_{i=1}^r \mathbf{u}_i^* \otimes \mathbf{v}_i^* \otimes \mathbf{z}_i \right\|_{\mathbb{F}} \leq \frac{11}{10} \lambda_{\max}^{*2/3} \|\mathbf{Z}\|_{\mathbb{F}}.
\end{aligned}$$

Combining these bounds, we reach

$$\gamma_1 \leq \frac{3}{2} \lambda_{\max}^{*4/3} \left(\|\mathbf{X}\|_{\mathbb{F}}^2 + \|\mathbf{Y}\|_{\mathbb{F}}^2 + \|\mathbf{Z}\|_{\mathbb{F}}^2 \right).$$

2. We now move on to γ_2 . Recall our assumption that $\|\mathbf{u}_i^*\|_2 = \|\mathbf{v}_i^*\|_2 = \|\mathbf{w}_i^*\|_2 = \lambda_i^{*1/3}$ for all $1 \leq i \leq r$. It follows from the Cauchy-Schwartz inequality that

$$\begin{aligned}
0 \leq \gamma_2 & \leq \frac{2}{3} \sum_{i=1}^r \lambda_i^{*2/3} \|\mathbf{x}_i\|^2 \|\mathbf{u}_i^*\|_2^2 + \frac{2}{3} \sum_{i=1}^r \lambda_i^{*2/3} \|\mathbf{y}_i\|^2 \|\mathbf{v}_i^*\|_2^2 + \frac{2}{3} \sum_{i=1}^r \lambda_i^{*2/3} \|\mathbf{z}_i\|^2 \|\mathbf{w}_i^*\|_2^2 \\
& \leq \frac{2}{3} \lambda_{\max}^{*4/3} \left(\|\mathbf{X}\|_{\mathbb{F}}^2 + \|\mathbf{Y}\|_{\mathbb{F}}^2 + \|\mathbf{Z}\|_{\mathbb{F}}^2 \right).
\end{aligned}$$

3. Turning to γ_3 , one can straightforwardly bound

$$|\gamma_3| \stackrel{(i)}{=} \frac{2}{3} \left| \sum_{i_1 \neq i_2} \langle \mathbf{x}_{i_1}, \mathbf{u}_{i_2}^* \rangle \langle \mathbf{v}_{i_1}^*, \mathbf{y}_{i_2} \rangle \langle \mathbf{w}_{i_1}^*, \mathbf{w}_{i_2}^* \rangle \right|$$

$$\begin{aligned}
&\stackrel{\text{(ii)}}{\leq} \max_{i_1 \neq i_2} |\langle \mathbf{w}_{i_1}^*, \mathbf{w}_{i_2}^* \rangle| \left(\sum_{i=1}^r \|\mathbf{x}_i\|_2 \|\mathbf{v}_i^*\|_2 \right) \left(\sum_{i=1}^r \|\mathbf{u}_i^*\|_2 \|\mathbf{y}_i\|_2 \right) \\
&\stackrel{\text{(iii)}}{\leq} \max_{i_1 \neq i_2} |\langle \mathbf{w}_{i_1}^*, \mathbf{w}_{i_2}^* \rangle| \| \mathbf{U}^* \|_{\text{F}} \| \mathbf{V}^* \|_{\text{F}} \| \mathbf{X} \|_{\text{F}} \| \mathbf{Y} \|_{\text{F}} \\
&\stackrel{\text{(iv)}}{\leq} r \sqrt{\frac{\mu}{d}} \lambda_{\max}^{*4/3} \| \mathbf{X} \|_{\text{F}} \| \mathbf{Y} \|_{\text{F}} \\
&\stackrel{\text{(v)}}{\ll} \lambda_{\max}^{*4/3} \| \mathbf{X} \|_{\text{F}} \| \mathbf{Y} \|_{\text{F}}.
\end{aligned}$$

Here, we have used the fact that $\|\mathbf{w}_i^*\|_2 = \lambda_i^{*1/3}$ in (i); the inequalities (ii) and (iii) arise from the Cauchy-Schwartz inequality; (iv) follows from (280c); (v) holds as long as $r \ll \sqrt{d_{\min}/\mu}$. In a similar manner, one can easily verify that

$$|\gamma_4| \ll \lambda_{\max}^{*4/3} \| \mathbf{X} \|_{\text{F}} \| \mathbf{Z} \|_{\text{F}} \quad \text{and} \quad |\gamma_5| \ll \lambda_{\max}^{*4/3} \| \mathbf{Y} \|_{\text{F}} \| \mathbf{Z} \|_{\text{F}}.$$

It then follows from the AM-GM inequality that

$$|\gamma_3| + |\gamma_4| + |\gamma_5| \ll \lambda_{\max}^{*4/3} \left(\| \mathbf{X} \|_{\text{F}}^2 + \| \mathbf{Y} \|_{\text{F}}^2 + \| \mathbf{Z} \|_{\text{F}}^2 \right).$$

4. Putting the above bounds together allows us to bound β_1 as follows

$$\begin{aligned}
\beta_1 &\leq \gamma_1 + \gamma_2 + |\gamma_3| + |\gamma_4| + |\gamma_5| \leq \frac{7}{2} \lambda_{\max}^{*4/3} \left(\| \mathbf{X} \|_{\text{F}}^2 + \| \mathbf{Y} \|_{\text{F}}^2 + \| \mathbf{Z} \|_{\text{F}}^2 \right); \\
\beta_1 &\geq \gamma_1 - |\gamma_3| - |\gamma_4| - |\gamma_5| \geq \frac{9}{10} \lambda_{\min}^{*4/3} \left(\| \mathbf{X} \|_{\text{F}}^2 + \| \mathbf{Y} \|_{\text{F}}^2 + \| \mathbf{Z} \|_{\text{F}}^2 \right).
\end{aligned}$$

Bounding β_2 To control β_2 , we note that β_2 involves two quantities: (1) the deviation of $\frac{1}{p} \left\| \mathcal{P}_{\Omega} \left(\sum_{i=1}^r \mathbf{x}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i + \mathbf{u}_i \otimes \mathbf{y}_i \otimes \mathbf{w}_i + \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{z}_i \right) \right\|_{\text{F}}^2$ from its expectation with respect to Ω ; (2) the distance between $\sum_{i=1}^r \mathbf{x}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i + \mathbf{u}_i \otimes \mathbf{y}_i \otimes \mathbf{w}_i + \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{z}_i$ and $\sum_{i=1}^r \mathbf{x}_i \otimes \mathbf{v}_i^* \otimes \mathbf{w}_i - \mathbf{u}_i^* \otimes \mathbf{y}_i \otimes \mathbf{w}_i^* - \mathbf{u}_i^* \otimes \mathbf{v}_i^* \otimes \mathbf{z}_i$. The first term can be shown to be exceedingly small under our sample size condition with the help of [YZ16, Lemma 5], while the second term is also guaranteed to be sufficiently small by the assumptions of the error of \mathbf{U} , \mathbf{V} and \mathbf{W} in (285). Therefore, one can apply an analogous argument for α_1 and α_2 in the proof of Lemma 5.1 in Appendix A to derive: with probability at least $1 - O(d_{\min}^{-10})$ one has

$$|\beta_2| \ll \lambda_{\min}^{*4/3} \left(\| \mathbf{X} \|_{\text{F}}^2 + \| \mathbf{Y} \|_{\text{F}}^2 + \| \mathbf{Z} \|_{\text{F}}^2 \right),$$

provided the sampling rate exceeds $p \gg \frac{\mu^2 r^2 d_{\max} \log^2 d_{\max}}{d_1 d_2 d_3}$.

Bounding β_3 By the definition of the operator norm, one can bound

$$\begin{aligned}
|\beta_3| &\lesssim \left\| p^{-1} \mathcal{P}_{\Omega} \left(\sum_{i=1}^r \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i - \mathbf{T}^* \right) \right\| \cdot \sum_{i=1}^r \left(\|\mathbf{w}_i\|_2 \|\mathbf{x}_i\|_2 \|\mathbf{y}_i\|_2 + \|\mathbf{v}_i\|_2 \|\mathbf{x}_i\|_2 \|\mathbf{z}_i\|_2 + \|\mathbf{u}_i\|_2 \|\mathbf{y}_i\|_2 \|\mathbf{z}_i\|_2 \right) \\
&\lesssim \lambda_{\max}^{2/3} \left\| \sum_{i=1}^r \mathbf{u}_i \otimes \mathbf{v}_i \otimes \mathbf{w}_i - \mathbf{T}^* \right\| \left\| p^{-1} \mathcal{P}_{\Omega} (\mathbf{1}_{d_1} \otimes \mathbf{1}_{d_2} \otimes \mathbf{1}_{d_3}) \right\| \left(\| \mathbf{X} \|_{\text{F}}^2 + \| \mathbf{Y} \|_{\text{F}}^2 + \| \mathbf{Z} \|_{\text{F}}^2 \right).
\end{aligned}$$

One can easily adapt the proof of Lemma D.2 to show that, with probability at least $1 - O(d_{\min}^{-10})$,

$$\left\| (p^{-1} \mathcal{P}_{\Omega} - \mathcal{I}) (\mathbf{1}_{d_1} \otimes \mathbf{1}_{d_2} \otimes \mathbf{1}_{d_3}) \right\| \lesssim \frac{\log^3 d_{\max}}{p} + \sqrt{\frac{d_{\max} \log^5 d_{\max}}{p}} \lesssim \sqrt{d_1 d_2 d_3} = \| \mathbf{1}_{d_1} \otimes \mathbf{1}_{d_2} \otimes \mathbf{1}_{d_3} \|$$

holds as long as $p \gg \max \left\{ \frac{\log^3 d_{\max}}{\sqrt{d_1 d_2 d_3}}, \frac{d_{\max} \log^5 d_{\max}}{d_1 d_2 d_3} \right\}$. We can then adapt the proof for bounding α_3 in the proof of Lemma 5.1 in Appendix A to derive that with probability at least $1 - O(d_{\min}^{-10})$,

$$|\beta_3| \ll \lambda_{\min}^{*4/3} \left(\| \mathbf{X} \|_{\text{F}}^2 + \| \mathbf{Y} \|_{\text{F}}^2 + \| \mathbf{Z} \|_{\text{F}}^2 \right).$$

Bounding β_4 It remains to control β_4 . By symmetry, it suffices to consider the following terms:

$$\begin{aligned}\gamma_1 &:= \sum_{i=1}^r (\alpha_i - \lambda_i^{*2/3}) (\langle \mathbf{x}_i, \mathbf{u}_i^* \rangle + \langle \mathbf{y}_i, \mathbf{v}_i^* \rangle)^2, \\ \gamma_2 &:= \sum_{i=1}^r \alpha_i (\langle \mathbf{x}_i, \Delta_i^U \rangle - \langle \mathbf{y}_i, \Delta_i^V \rangle)^2, \\ \gamma_3 &:= \sum_{i=1}^r \alpha_i (\langle \mathbf{x}_i, \mathbf{u}_i^* \rangle + \langle \mathbf{y}_i, \mathbf{v}_i^* \rangle) (\langle \mathbf{x}_i, \Delta_i^U \rangle - \langle \mathbf{y}_i, \Delta_i^V \rangle), \\ \gamma_4 &:= \sum_{i=1}^r \alpha_i (\|\mathbf{u}_i\|_2^2 - \|\mathbf{v}_i\|_2^2) (\|\mathbf{x}_i\|_2^2 - \|\mathbf{y}_i\|_2^2).\end{aligned}$$

In the following, we shall bound these four terms separately. By the assumptions that $\mathbf{S}^{(1)} = \mathbf{I}_r$, $\mathbf{S}^{(2)} = \mathbf{S}^{(3)} = -\mathbf{I}_r$, $\delta \ll 1$ and $\kappa \asymp 1$, one has

$$\|\Delta_i^U\|_2 \leq \|\mathbf{U} - \mathbf{U}^*\|_F \leq \delta \|\mathbf{U}^*\|_F \leq \delta \sqrt{r} \lambda_{\max}^{*1/3} \ll \lambda_{\min}^{*1/3} \leq \|\mathbf{u}_i^*\|_2; \quad (288a)$$

$$\|\Delta_i^V\|_2 \leq \|\mathbf{V} - \mathbf{V}^*\|_F \leq \delta \|\mathbf{V}^*\|_F \leq \delta \sqrt{r} \lambda_{\max}^{*1/3} \ll \lambda_{\min}^{*1/3} \leq \|\mathbf{v}_i^*\|_2; \quad (288b)$$

$$\|\Delta_i^W\|_2 \leq \|\mathbf{W} - \mathbf{W}^*\|_F \leq \delta \|\mathbf{W}^*\|_F \leq \delta \sqrt{r} \lambda_{\max}^{*1/3} \ll \lambda_{\min}^{*1/3} \leq \|\mathbf{w}_i^*\|_2. \quad (288c)$$

In addition, for each $1 \leq i \leq r$, one has

$$\frac{9}{10} \lambda_{\min}^{*2/3} \leq \lambda_i - |\alpha_i - \lambda_i^{*2/3}| \leq \alpha_i \leq \lambda_i + |\alpha_i - \lambda_i^{*2/3}| \leq \frac{11}{10} \lambda_{\max}^{*2/3} \quad (289)$$

by virtue of the condition that $|\alpha_i - \lambda_i^{*2/3}| \ll \lambda_{\min}^{*2/3}$.

1. Let us start with γ_1 . By the condition that $\max_i |\alpha_i - \lambda_i^{*2/3}| \ll \lambda_{\min}^{*2/3}$, it is easily seen that

$$\begin{aligned}|\gamma_1| &\lesssim \sum_{i=1}^r |\alpha_i - \lambda_i^{*2/3}| (\langle \mathbf{x}_i, \mathbf{u}_i^* \rangle^2 + \langle \mathbf{y}_i, \mathbf{v}_i^* \rangle^2) \leq \max_i |\alpha_i - \lambda_i^{*2/3}| \sum_{i=1}^r (\|\mathbf{x}_i\|_2^2 \|\mathbf{u}_i^*\|_2^2 + \|\mathbf{y}_i\|_2^2 \|\mathbf{v}_i^*\|_2^2) \\ &\leq \max_i |\alpha_i - \lambda_i^{*2/3}| (\|\mathbf{u}_i^*\|_2^2 + \|\mathbf{v}_i^*\|_2^2) (\|\mathbf{X}\|_F^2 + \|\mathbf{Y}\|_F^2) \ll \lambda_{\min}^{*4/3} (\|\mathbf{X}\|_F^2 + \|\mathbf{Y}\|_F^2),\end{aligned}$$

with the proviso that $\kappa \asymp 1$.

2. As for γ_2 , one can bound

$$\begin{aligned}|\gamma_2| &\lesssim \sum_{i=1}^r \alpha_i (\langle \mathbf{x}_i, \Delta_i^U \rangle^2 - \langle \mathbf{y}_i, \Delta_i^V \rangle^2) \leq \max_i \alpha_i \sum_{i=1}^r (\|\mathbf{x}_i\|_2^2 \|\Delta_i^U\|_2^2 + \|\mathbf{y}_i\|_2^2 \|\Delta_i^V\|_2^2) \\ &\leq \max_i \alpha_i (\|\Delta_i^U\|_2^2 + \|\Delta_i^V\|_2^2) (\|\mathbf{X}\|_F^2 + \|\mathbf{Y}\|_F^2) \ll \lambda_{\min}^{*4/3} (\|\mathbf{X}\|_F^2 + \|\mathbf{Y}\|_F^2).\end{aligned}$$

Here, we use (288a), (288b), (289) and $\kappa \asymp 1$ in the last step.

3. Turning to γ_3 , one can develop a similar bound as follows

$$\begin{aligned}|\gamma_3| &\leq \sqrt{\sum_{i=1}^r \alpha_i (\langle \mathbf{x}_i, \mathbf{u}_i^* \rangle + \langle \mathbf{y}_i, \mathbf{v}_i^* \rangle)^2} \cdot \sqrt{\sum_{i=1}^r \alpha_i (\langle \mathbf{x}_i, \Delta_i^U \rangle - \langle \mathbf{y}_i, \Delta_i^V \rangle)^2} \\ &\lesssim \max_i \alpha_i \sqrt{(\|\mathbf{u}_i^*\|_2^2 + \|\mathbf{v}_i^*\|_2^2) (\|\mathbf{X}\|_F^2 + \|\mathbf{Y}\|_F^2)} \cdot \sqrt{(\|\Delta_i^U\|_2^2 + \|\Delta_i^V\|_2^2) (\|\mathbf{X}\|_F^2 + \|\mathbf{Y}\|_F^2)} \\ &\ll \lambda_{\min}^{*4/3} (\|\mathbf{X}\|_F^2 + \|\mathbf{Y}\|_F^2).\end{aligned}$$

4. With regards to γ_4 , we can expand

$$\begin{aligned}
|\gamma_4| &= \sum_{i=1}^r \alpha_i (\|\mathbf{u}_i^* + \Delta_i^U\|_2^2 - \|\mathbf{v}_i^* + \Delta_i^V\|_2^2) (\|\mathbf{x}_i\|_2^2 - \|\mathbf{y}_i\|_2^2) \\
&= \sum_{i=1}^r \alpha_i (\|\mathbf{u}_i^*\|_2^2 + 2\langle \mathbf{u}_i^*, \Delta_i^U \rangle + \|\Delta_i^U\|_2^2 - \|\mathbf{v}_i^*\|_2^2 + 2\langle \mathbf{v}_i^*, \Delta_i^V \rangle - \|\Delta_i^V\|_2^2) (\|\mathbf{x}_i\|_2^2 - \|\mathbf{y}_i\|_2^2) \\
&= \sum_{i=1}^r \alpha_i (2\langle \mathbf{u}_i^*, \Delta_i^U \rangle + \|\Delta_i^U\|_2^2 + 2\langle \mathbf{v}_i^*, \Delta_i^V \rangle - \|\Delta_i^V\|_2^2) (\|\mathbf{x}_i\|_2^2 - \|\mathbf{y}_i\|_2^2),
\end{aligned}$$

where the last step follows from the assumption $\|\mathbf{u}_i^*\|_2 = \|\mathbf{v}_i^*\|_2 = \|\mathbf{w}_i^*\|_2$. It follows that

$$\begin{aligned}
\left| \sum_{i=1}^r \alpha_i (\|\mathbf{x}_i\|_2^2 - \|\mathbf{y}_i\|_2^2) (\|\mathbf{x}_i\|_2^2 - \|\mathbf{y}_i\|_2^2) \right| &\lesssim \sum_{i=1}^r \alpha_i |\langle \mathbf{u}_i^*, \Delta_i^U \rangle + \|\Delta_i^U\|_2^2 + \langle \mathbf{v}_i^*, \Delta_i^V \rangle + \|\Delta_i^V\|_2^2| (\|\mathbf{x}_i\|_2^2 + \|\mathbf{y}_i\|_2^2) \\
&\leq \max_i \alpha_i (\|\mathbf{u}_i^*\|_2 \|\Delta_i^U\|_2 + \|\mathbf{v}_i^*\|_2 \|\Delta_i^V\|_2) \sum_{i=1}^r (\|\mathbf{x}_i\|_2^2 + \|\mathbf{y}_i\|_2^2) \\
&\ll \lambda_{\min}^{*4/3} \left(\|\mathbf{X}\|_{\text{F}}^2 + \|\mathbf{Y}\|_{\text{F}}^2 \right),
\end{aligned}$$

where we have used the conditions that $\max_i \alpha_i \lesssim \lambda_{\max}^{*2/3}$ and $\kappa \asymp 1$.

5. It is not hard to check that similar bounds also hold for the remaining terms in β_4 . As a result, we arrive at

$$|\beta_4| \ll \lambda_{\min}^{*4/3} \left(\|\mathbf{X}\|_{\text{F}}^2 + \|\mathbf{Y}\|_{\text{F}}^2 + \|\mathbf{Z}\|_{\text{F}}^2 \right).$$

Combining the previous bounds on $\beta_1, \beta_2, \beta_3$ and β_4 Putting the above estimates together, we conclude that with probability at least $1 - O(d_{\min}^{-10})$, one has

$$\begin{aligned}
\text{vec} \left(\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \\ \mathbf{Z} \end{bmatrix} \right)^\top \nabla^2 g_{\text{clean}}(\mathbf{U}, \mathbf{V}, \mathbf{W}) \text{vec} \left(\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \\ \mathbf{Z} \end{bmatrix} \right) &\geq \beta_1 - |\beta_2| - |\beta_3| - |\beta_4| \geq \frac{1}{2} \lambda_{\min}^{*4/3} \left(\|\mathbf{X}\|_{\text{F}}^2 + \|\mathbf{Y}\|_{\text{F}}^2 + \|\mathbf{Z}\|_{\text{F}}^2 \right) \\
\text{vec} \left(\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \\ \mathbf{Z} \end{bmatrix} \right)^\top \nabla^2 g_{\text{clean}}(\mathbf{U}, \mathbf{V}, \mathbf{W}) \text{vec} \left(\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \\ \mathbf{Z} \end{bmatrix} \right) &\leq \beta_1 + |\beta_2| + |\beta_3| + |\beta_4| \leq 4\lambda_{\max}^{*4/3} \left(\|\mathbf{X}\|_{\text{F}}^2 + \|\mathbf{Y}\|_{\text{F}}^2 + \|\mathbf{Z}\|_{\text{F}}^2 \right)
\end{aligned}$$

as claimed.

References

- [AFWZ17] E. Abbe, J. Fan, K. Wang, and Y. Zhong. Entrywise eigenvector analysis of random matrices with low expected rank. *arXiv:1709.09565*, 2017.
- [AGH⁺14] A. Anandkumar, R. Ge, D. Hsu, S. M. Kakade, and M. Telgarsky. Tensor decompositions for learning latent variable models. *The Journal of Machine Learning Research*, 15(1):2773–2832, 2014.
- [AGJ14] A. Anandkumar, R. Ge, and M. Janzamin. Guaranteed non-orthogonal tensor decomposition via alternating rank-1 updates. *arXiv preprint arXiv:1402.5180*, 2014.
- [AGJ15] A. Anandkumar, R. Ge, and M. Janzamin. Learning overcomplete latent variable models through tensor methods. In *Proceedings of the Conference on Learning Theory*, pages 36–112, 2015.

- [AGJ17] A. Anandkumar, R. Ge, and M. Janzamin. Analyzing tensor power method dynamics in over-complete regime. *The Journal of Machine Learning Research*, 18(1):752–791, 2017.
- [AW17] M. Ashraphijuo and X. Wang. Fundamental conditions for low-CP-rank tensor completion. *The Journal of Machine Learning Research*, 18(1):2116–2145, 2017.
- [BM16] B. Barak and A. Moitra. Noisy tensor completion via the sum-of-squares hierarchy. In *Conference on Learning Theory*, pages 417–445, 2016.
- [Bub15] S. Bubeck. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- [CC14] Y. Chen and Y. Chi. Robust spectral compressed sensing via structured matrix completion. *IEEE Transactions on Information Theory*, 60(10):6576 – 6601, 2014.
- [CC17] Y. Chen and E. J. Candès. Solving random quadratic systems of equations is nearly as easy as solving linear systems. *Comm. Pure Appl. Math.*, 70(5):822–883, 2017.
- [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, July 2018.
- [CCF⁺19] Y. Chen, Y. Chi, J. Fan, C. Ma, and Y. Yan. Noisy matrix completion: Understanding statistical guarantees for convex relaxation via nonconvex optimization. *arXiv preprint arXiv:1902.07698*, accepted to *SIAM Journal on Optimization*, 2019.
- [CCFM19] Y. Chen, Y. Chi, J. Fan, and C. Ma. Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval. *Mathematical Programming*, 176(1-2):5–37, July 2019.
- [CFMW19] Y. Chen, J. Fan, C. Ma, and K. Wang. Spectral method and regularized MLE are both optimal for top- K ranking. *Annals of Statistics*, 47(4):2204–2235, August 2019.
- [CFMY19] Y. Chen, J. Fan, C. Ma, and Y. Yan. Inference and uncertainty quantification for noisy matrix completion. *arXiv:1906.04159*, accepted to the *Proceedings of the National Academy of Sciences (PNAS)*, 2019.
- [CKK⁺97] C. A. Cocosco, V. Kollokian, R. K.-S. Kwan, G. B. Pike, and A. C. Evans. Brainweb: Online interface to a 3d mri simulated brain database. In *NeuroImage*. Citeseer, 1997.
- [CLC19] Y. Chi, Y. Lu, and Y. Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing*, 67(20):5239–5269, 2019.
- [CLC⁺20] C. Cai, G. Li, Y. Chi, H. V. Poor, and Y. Chen. Subspace estimation from unbalanced and incomplete data matrices: $\ell_{2,\infty}$ statistical guarantees. *The Annals of Statistics*, to appear, 2020.
- [CLL19] J. Chen, D. Liu, and X. Li. Nonconvex rectangular matrix completion via gradient descent without $\ell_{2,\infty}$ regularization. *arXiv:1901.06116v1*, 2019.
- [CLPC19] C. Cai, G. Li, H. V. Poor, and Y. Chen. Nonconvex low-rank tensor completion from noisy data. In *Advances in Neural Information Processing Systems*, pages 1863–1874, 2019.
- [CLS15] E. J. Candès, X. Li, and M. Soltanolkotabi. Phase retrieval via wirtinger flow: Theory and algorithms. *IEEE Transactions on Information Theory*, 61(4):1985–2007, 2015.
- [CR09] E. Candès and B. Recht. Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9(6):717–772, April 2009.
- [CRY19] H. Chen, G. Raskutti, and M. Yuan. Non-convex projected gradient descent for generalized low-rank tensor regression. *The Journal of Machine Learning Research*, 20(1):172–208, 2019.

- [CW15] Y. Chen and M. J. Wainwright. Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees. *arXiv:1509.03025*, 2015.
- [CZA⁺17] J. Y. Cheng, T. Zhang, M. T. Alley, M. Uecker, M. Lustig, J. M. Pauly, and S. S. Vasanaawala. Comprehensive multi-dimensional MRI for the simultaneous assessment of cardiopulmonary anatomy and physiology. *Scientific reports*, 7(1):5330, 2017.
- [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.
- [DC18] L. Ding and Y. Chen. The leave-one-out approach for matrix completion: Primal and dual analysis. *arXiv preprint arXiv:1803.07554*, 2018.
- [DR16] M. A. Davenport and J. Romberg. An overview of low-rank matrix recovery from incomplete observations. *IEEE Journal of Selected Topics in Signal Processing*, 10(4):608–622, 2016.
- [EAHK13] G. Ely, S. Aeron, N. Hao, and M. E. Kilmer. 5D and 4D pre-stack seismic data completion using tensor nuclear norm (TNN). In *SEG Technical Program Expanded Abstracts 2013*, pages 3639–3644. Society of Exploration Geophysicists, 2013.
- [EK15] N. El Karoui. On the impact of predictor geometry on the performance on high-dimensional ridge-regularized generalized robust regression estimators. *Probability Theory and Related Fields*, pages 1–81, 2015.
- [FL19] V. F. Farias and A. A. Li. Learning preferences with side information. *Management Science*, 65(7):3131–3149, 2019.
- [GBW18] D. Gilboa, S. Buchanan, and J. Wright. Efficient dictionary learning with gradient descent. *arXiv preprint arXiv:1809.10313*, 2018.
- [GHJY15] R. Ge, F. Huang, C. Jin, and Y. Yuan. Escaping from saddle points online stochastic gradient for tensor decomposition. In *Conference on Learning Theory*, pages 797–842, 2015.
- [GM17] R. Ge and T. Ma. On the optimization landscape of tensor decompositions. *arXiv preprint arXiv:1706.05598*, 2017.
- [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.
- [Gro11] D. Gross. Recovering low-rank matrices from few coefficients in any basis. *IEEE Transactions on Information Theory*, 57(3):1548–1566, March 2011.
- [GRY11] S. Gandy, B. Recht, and I. Yamada. Tensor completion and low- n -rank tensor recovery via convex optimization. *Inverse Problems*, 27(2):025010, 2011.
- [HL13] C. J. Hillar and L.-H. Lim. Most tensor problems are np-hard. *Journal of the ACM (JACM)*, 60(6):45, 2013.
- [HMGW15] B. Huang, C. Mu, D. Goldfarb, and J. Wright. Provable models for robust low-rank tensor completion. *Pacific Journal of Optimization*, 11(2):339–364, 2015.
- [HSS16] S. B. Hopkins, T. Schramm, J. Shi, and D. Steurer. Fast spectral algorithms from sum-of-squares proofs: tensor decomposition and planted sparse vectors. In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pages 178–191. ACM, 2016.
- [HWW⁺19] B. Hao, B. Wang, P. Wang, J. Zhang, J. Yang, and W. W. Sun. Sparse tensor additive regression. *arXiv preprint arXiv:1904.00479*, 2019.
- [HZC20] B. Hao, A. Zhang, and G. Cheng. Sparse and low-rank tensor estimation via cubic sketchings. *IEEE Transactions on Information Theory*, 2020.

- [JHZ⁺16] T.-Y. Ji, T.-Z. Huang, X.-L. Zhao, T.-H. Ma, and G. Liu. Tensor completion using total variation and low-rank matrix factorization. *Information Sciences*, 326:243–257, 2016.
- [JO14] P. Jain and S. Oh. Provable tensor factorization with missing data. In *Advances in Neural Information Processing Systems*, pages 1431–1439, 2014.
- [KB09] T. G. Kolda and B. W. Bader. Tensor decompositions and applications. *SIAM review*, 51(3):455–500, 2009.
- [KBHH13] M. E. Kilmer, K. Braman, N. Hao, and R. C. Hoover. Third-order tensors as operators on matrices: A theoretical and computational framework with applications in imaging. *SIAM Journal on Matrix Analysis and Applications*, 34(1):148–172, 2013.
- [KM16] H. Kasai and B. Mishra. Low-rank tensor completion: a riemannian manifold preconditioning approach. In *International Conference on Machine Learning*, pages 1012–1021, 2016.
- [KMO10a] R. H. Keshavan, A. Montanari, and S. Oh. Matrix completion from a few entries. *IEEE transactions on information theory*, 56(6):2980–2998, 2010.
- [KMO10b] R. H. Keshavan, A. Montanari, and S. Oh. Matrix completion from noisy entries. *J. Mach. Learn. Res.*, 11:2057–2078, 2010.
- [KOKC13] H.-J. Kim, E. Ollila, V. Koivunen, and C. Croux. Robust and sparse estimation of tensor decompositions. In *2013 IEEE Global Conference on Signal and Information Processing*, pages 965–968. IEEE, 2013.
- [Kol01] T. G. Kolda. Orthogonal tensor decompositions. *SIAM Journal on Matrix Analysis and Applications*, 23(1):243–255, 2001.
- [Kol11] V. Koltchinskii. *Oracle inequalities in empirical risk minimization and sparse recovery problems*, volume 2033 of *Lecture Notes in Mathematics*. 2011.
- [KS13] A. Krishnamurthy and A. Singh. Low-rank matrix and tensor completion via adaptive sampling. In *Advances in Neural Information Processing Systems*, pages 836–844, 2013.
- [KSS13] N. Kreimer, A. Stanton, and M. D. Sacchi. Tensor completion based on nuclear norm minimization for 5d seismic data reconstruction. *Geophysics*, 78(6):V273–V284, 2013.
- [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.
- [LMWY13] 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, 2013.
- [Lou14] K. Lounici. High-dimensional covariance matrix estimation with missing observations. *Bernoulli*, 20(3):1029–1058, 2014.
- [LSC⁺14] Y. Liu, F. Shang, H. Cheng, J. Cheng, and H. Tong. Factor matrix trace norm minimization for low-rank tensor completion. In *Proceedings of the 2014 SIAM International Conference on Data Mining*, pages 866–874. SIAM, 2014.
- [LT17] Q. Li and G. Tang. Convex and nonconvex geometries of symmetric tensor factorization. In *2017 51st Asilomar Conference on Signals, Systems, and Computers*, pages 305–309. IEEE, 2017.
- [LYX17] X. Li, Y. Ye, and X. Xu. Low-rank tensor completion with total variation for visual data inpainting. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.

- [LZT19] Q. Li, Z. Zhu, and G. Tang. The non-convex geometry of low-rank matrix optimization. *Information and Inference: A Journal of the IMA*, 8(1):51–96, 2019.
- [MHWG14] C. Mu, B. Huang, J. Wright, and D. Goldfarb. Square deal: Lower bounds and improved relaxations for tensor recovery. In *International conference on machine learning*, pages 73–81, 2014.
- [MP20] V. V. Mišić and G. Perakis. Data analytics in operations management: A review. *Manufacturing & Service Operations Management*, 22(1):158–169, 2020.
- [MS18] A. Montanari and N. Sun. Spectral algorithms for tensor completion. *Communications on Pure and Applied Mathematics*, 71(11):2381–2425, 2018.
- [MWCC17] C. Ma, K. Wang, Y. Chi, and Y. Chen. Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion and blind deconvolution. *arXiv preprint arXiv:1711.10467*, accepted to *Foundations of Computational Mathematics*, 2017.
- [NDT15] N. H. Nguyen, P. Drineas, and T. D. Tran. Tensor sparsification via a bound on the spectral norm of random tensors. *Information and Inference: A Journal of the IMA*, 4(3):195–229, 2015.
- [Paw19] C. Pawlowski. *Machine learning for problems with missing and uncertain data with applications to personalized medicine*. PhD thesis, Massachusetts Institute of Technology, 2019.
- [PS17] A. Potechin and D. Steurer. Exact tensor completion with sum-of-squares. In *Conference on Learning Theory*, pages 1619–1673, 2017.
- [PW19] A. Pananjady and M. J. Wainwright. Value function estimation in markov reward processes: Instance-dependent ℓ_∞ -bounds for policy evaluation. *arXiv preprint arXiv:1909.08749*, 2019.
- [RM14] E. Richard and A. Montanari. A statistical model for tensor PCA. In *Advances in Neural Information Processing Systems*, pages 2897–2905, 2014.
- [RPP13] B. Romera-Paredes and M. Pontil. A new convex relaxation for tensor completion. In *Advances in Neural Information Processing Systems*, pages 2967–2975, 2013.
- [RSS17] H. Rauhut, R. Schneider, and Ž. Stojanac. Low rank tensor recovery via iterative hard thresholding. *Linear Algebra and its Applications*, 523:220–262, 2017.
- [SDLF⁺17] N. D. Sidiropoulos, L. De Lathauwer, X. Fu, K. Huang, E. E. Papalexakis, and C. Faloutsos. Tensor decomposition for signal processing and machine learning. *IEEE Transactions on Signal Processing*, 65(13):3551–3582, 2017.
- [SHKM14] O. Semerci, N. Hao, M. E. Kilmer, and E. L. Miller. Tensor-based formulation and nuclear norm regularization for multienergy computed tomography. *IEEE Transactions on Image Processing*, 23(4):1678–1693, 2014.
- [SL16] R. Sun and Z.-Q. Luo. Guaranteed matrix completion via non-convex factorization. *IEEE Transactions on Information Theory*, 62(11):6535–6579, 2016.
- [SLLC17] W. W. Sun, J. Lu, H. Liu, and G. Cheng. Provable sparse tensor decomposition. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(3):899–916, 2017.
- [SN16] S. R. Soroushmehr and K. Najarian. Transforming big data into computational models for personalized medicine and health care. *Dialogues in clinical neuroscience*, 18(3):339, 2016.
- [Ste16] M. Steinlechner. Riemannian optimization for high-dimensional tensor completion. *SIAM Journal on Scientific Computing*, 38(5):S461–S484, 2016.
- [SY19] D. Shah and C. L. Yu. Iterative collaborative filtering for sparse noisy tensor estimation. *arXiv preprint arXiv:1908.01241*, 2019.

- [TBS⁺16] S. Tu, R. Boczar, M. Simchowitz, M. Soltanolkotabi, and B. Recht. Low-rank solutions of linear matrix equations via procrustes flow. In *International Conference on Machine Learning*, pages 964–973, 2016.
- [THK10] R. Tomioka, K. Hayashi, and H. Kashima. Estimation of low-rank tensors via convex optimization. *arXiv preprint arXiv:1010.0789*, 2010.
- [TS15] G. Tang and P. Shah. Guaranteed tensor decomposition: A moment approach. In *International Conference on Machine Learning*, pages 1491–1500, 2015.
- [TV19] Y. S. Tan and R. Vershynin. Online stochastic gradient descent with arbitrary initialization solves non-smooth, non-convex phase retrieval. *arXiv preprint arXiv:1910.12837*, 2019.
- [Ver10] R. Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- [Ver18] R. Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge University Press, 2018.
- [WAA16] W. Wang, V. Aggarwal, and S. Aeron. Tensor completion by alternating minimization under the tensor train (tt) model. *arXiv preprint arXiv:1609.05587*, 2016.
- [XHYS15] Y. Xu, R. Hao, W. Yin, and Z. Su. Parallel matrix factorization for low-rank tensor completion. *Inverse Problems & Imaging*, 9(2):601–624, 2015.
- [XY13] Y. Xu and W. Yin. A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion. *SIAM Journal on imaging sciences*, 6(3):1758–1789, 2013.
- [XY17] D. Xia and M. Yuan. On polynomial time methods for exact low rank tensor completion. *arXiv preprint arXiv:1702.06980*, 2017.
- [XYZ17] D. Xia, M. Yuan, and C.-H. Zhang. Statistically optimal and computationally efficient low rank tensor completion from noisy entries. *arXiv preprint arXiv:1711.04934*, 2017.
- [Yao18] Q. Yao. Scalable tensor completion with nonconvex regularization. *arXiv preprint arXiv:1807.08725*, 2018.
- [YLW⁺17] J. Ying, H. Lu, Q. Wei, J.-F. Cai, D. Guo, J. Wu, Z. Chen, and X. Qu. Hankel matrix nuclear norm regularized tensor completion for n -dimensional exponential signals. *IEEE Transactions on Signal Processing*, 65(14):3702–3717, 2017.
- [YPCC16] X. Yi, D. Park, Y. Chen, and C. Caramanis. Fast algorithms for robust PCA via gradient descent. In *NIPS*, 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.
- [YZ17] M. Yuan and C.-H. Zhang. Incoherent tensor norms and their applications in higher order tensor completion. *IEEE Transactions on Information Theory*, 63(10):6753–6766, 2017.
- [ZA17] Z. Zhang and S. Aeron. Exact tensor completion using t-svd. *IEEE Trans. Signal Processing*, 65(6):1511–1526, 2017.
- [ZB18] Y. Zhong and N. Boumal. Near-optimal bounds for phase synchronization. *SIAM Journal on Optimization*, 28(2):989–1016, 2018.
- [Zha19] A. Zhang. Cross: Efficient low-rank tensor completion. *The Annals of Statistics*, 47(2):936–964, 2019.

- [ZKOM18] A. M. Zoubir, V. Koivunen, E. Ollila, and M. Muma. *Robust statistics for signal processing*. Cambridge University Press, 2018.
- [ZL16] Q. Zheng and J. Lafferty. Convergence analysis for rectangular matrix completion using Burer-Monteiro factorization and gradient descent. *arXiv:1605.07051*, 2016.
- [ZX18] A. Zhang and D. Xia. Tensor SVD: Statistical and computational limits. *IEEE Transactions on Information Theory*, 64(11):7311–7338, 2018.