

# Learning Mixtures of Low-Rank Models

Yanxi Chen\*    Cong Ma†    H. Vincent Poor\*    Yuxin Chen\*

September 23, 2020

## Abstract

We study the problem of learning mixtures of low-rank models, i.e. reconstructing multiple low-rank matrices from unlabelled linear measurements of each. This problem enriches two widely studied settings — low-rank matrix sensing and mixed linear regression — by bringing latent variables (i.e. unknown labels) and structural priors (i.e. low-rank structures) into consideration. To cope with the non-convexity issues arising from unlabelled heterogeneous data and low-complexity structure, we develop a three-stage meta-algorithm that is guaranteed to recover the unknown matrices with near-optimal sample and computational complexities under Gaussian designs. In addition, the proposed algorithm is provably stable against random noise. We complement the theoretical studies with empirical evidence that confirms the efficacy of our algorithm.

**Keywords:** matrix sensing, latent variable models, heterogeneous data, mixed linear regression, non-convex optimization, meta-learning

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Main contributions	3
1.2	Notation	3
<b>2</b>	<b>Algorithm</b>	<b>4</b>
2.1	Stage 1: subspace estimation via a spectral method	4
2.2	Stage 2: initialization via low-dimensional mixed linear regression	4
2.3	Stage 3: local refinement via scaled truncated gradient descent (ScaledTGD)	6
2.4	The full algorithm	6
<b>3</b>	<b>Main results</b>	<b>7</b>
3.1	Models and assumptions	7
3.2	Theoretical guarantees	8
3.3	Numerical experiments	9
<b>4</b>	<b>Prior work</b>	<b>10</b>
<b>5</b>	<b>Analysis</b>	<b>11</b>
<b>6</b>	<b>Discussion</b>	<b>12</b>
<b>A</b>	<b>The tensor method for mixed linear regression</b>	<b>13</b>

---

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

†Department of Electrical Engineering and Computer Science, University of California Berkeley, Berkeley, CA 94720, USA; email: congm@berkeley.edu

<b>B Proofs for Section 5</b>	<b>14</b>
B.1 Proof of Theorem 3	14
B.2 Proof of Theorem 4	16
B.3 Proof of Proposition 1	19
B.4 Proof of Theorem 5	19
<b>C Technical lemmas</b>	<b>24</b>
C.1 Proof of Lemma 1	26
C.2 Proof of Lemma 2	30
C.3 Proof of Lemma 3	32
<b>D Estimating unknown parameters in Algorithm 4</b>	<b>33</b>

## 1 Introduction

This paper explores a mixture of low-rank models with latent variables, which seeks to reconstruct a couple of low-rank matrices  $\mathbf{M}_k^* \in \mathbb{R}^{n_1 \times n_2}$  ( $1 \leq k \leq K$ ) from *unlabeled* linear measurements of each. More specifically, what we have available is a collection of  $N$  linear measurements  $\{y_i\}_{1 \leq i \leq N}$  taking the following form:

$$y_i = \begin{cases} \langle \mathbf{A}_i, \mathbf{M}_1^* \rangle, & \text{if } i \in \Omega_1^*, \\ \dots & \dots \\ \langle \mathbf{A}_i, \mathbf{M}_K^* \rangle, & \text{if } i \in \Omega_K^*, \end{cases} \quad (1)$$

where  $\{\mathbf{A}_i\}_{1 \leq i \leq N}$  are the sampling/design matrices,  $\langle \cdot, \cdot \rangle$  denotes the matrix inner product, and  $\{\Omega_k^*\}_{1 \leq k \leq K}$  represents an unknown partition of the index set  $\{1, \dots, N\}$ . The aim is to design an algorithm that is guaranteed to recover  $\{\mathbf{M}_k^*\}$  efficiently and faithfully, despite the absence of knowledge of  $\{\Omega_k^*\}_{1 \leq k \leq K}$ .

This problem of learning mixtures of low-rank models enriches two widely studied settings: (1) it generalizes classical low-rank matrix recovery [RFP10, CLC19] by incorporating heterogeneous data and latent variables (i.e. the labels indicating which low-rank matrices are being measured), and (2) it expands the studies of mixed linear regression [QR78, YCS14] by integrating low-complexity structural priors (i.e. low-rank structures). In addition to the prior work [YC15] that has studied this setting, we single out two broader scenarios that bear relevance to and motivate the investigation of mixtures of low-rank models.

- *Mixed matrix completion.* If each measurement  $y_i$  only reveals a single entry of one of the unknown matrices  $\{\mathbf{M}_k^*\}$ , then the problem is commonly referred to as mixed matrix completion (namely, completing several low-rank matrices from a mixture of unlabeled observations of their entries) [PA18]. One motivating application arises from computer vision, where several problems like joint shape matching can be posed as structured matrix completion [CGH14, CC18a]. When the objects to be matched exhibit certain geometric symmetry, there might exist multiple plausible maps (and hence multiple ground-truth matrices), and the provided observations might become intrinsically unlabeled due to symmetric ambiguities [SLHH18]. Other applications include network topology inference and metagenomics given mixed DNA samples; see [PA18] for details.
- *Multi-task learning and meta-learning.* The model (1) can be viewed as an instance of multi-task learning or meta-learning [Bax00, MPRP16, KSS<sup>+</sup>20], where the tasks follow a discrete prior distribution supported on a set of  $K$  meta parameters, and each training data point  $(\mathbf{A}_i, y_i)$  is a realization of one task that comes with a single sample. While it is typically assumed in meta-learning that even light tasks have more than one samples, understanding this single-sample model is essential towards tackling more general settings. Additionally, in comparison to meta-learning for mixed linear regression [KSS<sup>+</sup>20, KSKO20], the model (1) imposes further structural prior on the unknown meta parameters, thereby allowing for potential reduction of sample complexities.

The challenge for learning mixtures of low-rank models primarily stems from the non-convexity issues. While the low-rank structure alone already leads to non-convex optimization landscapes, the presence of heterogeneous data and discrete hidden variables further complicates matters significantly.

## 1.1 Main contributions

This paper takes a step towards learning mixtures of low-rank models, focusing on the tractable Gaussian design where the  $\mathbf{A}_i$ 's have i.i.d. Gaussian entries; in light of this, we shall also call the problem *mixed matrix sensing*, to be consistent with the terminology used in recent literature [BNS16, CLC19]. In particular, we propose a meta-algorithm comprising the following three stages:

1. Estimate the joint column and row spaces of  $\{\mathbf{M}_k^*\}_{1 \leq k \leq K}$ ;
2. Transform mixed matrix sensing into low-dimensional mixed linear regression using the above subspace estimates, and invoke a mixed linear regression solver to obtain initial estimates of  $\{\mathbf{M}_k^*\}_{1 \leq k \leq K}$ ;
3. Successively refine the estimates via a non-convex low-rank matrix factorization algorithm (more specifically, an algorithm called *scaled truncated gradient descent* to be described in Algorithm 3).

The details of each stage will be spelled out and elucidated in Section 2.

Encouragingly, the proposed algorithm is guaranteed to succeed under mild conditions (to be specified in Section 3.1). Informally, our contributions are three-fold.

- *Exact recovery in the noiseless case.* In the absence of noise, our algorithm enables exact recovery of  $\{\mathbf{M}_k^*\}$  modulo global permutation. The sample complexity required to achieve this scales linearly (up to some log factor) in the dimension  $\max\{n_1, n_2\}$  and polynomially in other salient parameters.
- *Stability vis-à-vis random noise.* The proposed algorithm is provably stable against Gaussian noise, in the sense that the estimation accuracy degrades gracefully as the signal-to-noise-ratio decreases.
- *Computational efficiency.* When the number  $K$  of components and the maximum rank of the unknown matrices are both constants, the computational cost of our algorithm scales nearly linearly in  $Nn_1n_2$  with  $N$  the number of samples — this is proportional to the time taken to read all design matrices.

The precise theorem statements are postponed to Section 3. Empirical evidence will also be provided in Section 3 to corroborate the efficacy of our algorithm.

## 1.2 Notation

Before we proceed, let us collect some notation that will be frequently used. Throughout this paper, we reserve boldfaced symbols for vectors (lower case) and matrices (upper case). For a vector  $\mathbf{x}$ ,  $\|\mathbf{x}\|_2$  denotes its  $\ell_2$  norm. For a matrix  $\mathbf{X}$ ,  $\|\mathbf{X}\|$  (resp.  $\|\mathbf{X}\|_F$ ) denotes its spectral (resp. Frobenius) norm,  $\sigma_k(\mathbf{X})$  denotes its  $k$ -th largest singular value, and  $\text{col}\{\mathbf{X}\}$  (resp.  $\text{row}\{\mathbf{X}\}$ ) denotes its column (resp. row) space. If  $\mathbf{U}$  is a matrix with orthonormal columns, we also use the same notation  $\mathbf{U}$  to represent its column space, and vice versa. For any matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n_1 \times n_2}$ , let  $\langle \mathbf{A}, \mathbf{B} \rangle := \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} A_{ij} B_{ij}$  stand for the matrix inner product.  $\mathbf{I}_n$  represents the  $n \times n$  identity matrix.  $\text{vec}(\cdot)$  denotes vectorization of a matrix, and  $\text{mat}(\cdot)$  denotes the inverse operation (the corresponding matrix dimensions should often be clear from the context).

We use both  $a_n \lesssim b_n$  and  $a_n = O(b_n)$  to indicate that  $a_n \leq C_0 b_n$  for some universal constant  $C_0 > 0$ ; in addition,  $a_n \gtrsim b_n$  is equivalent to  $b_n \lesssim a_n$ , and  $a_n \asymp b_n$  means both  $a_n \lesssim b_n$  and  $b_n \lesssim a_n$  hold true. Finally,  $a_n = o(b_n)$  means that  $a_n/b_n \rightarrow 0$  as  $n \rightarrow \infty$ .

For a finite set  $\Omega$ , we denote by  $|\Omega|$  its cardinality. For a number  $\alpha \in [0, 1]$  and a random variable  $X$  following some distribution on  $\mathbb{R}$ , we let  $Q_\alpha(X)$  denote the  $\alpha$ -quantile function, namely

$$Q_\alpha(X) := \inf \{t \in \mathbb{R} : \mathbb{P}(X \leq t) \geq \alpha\}. \quad (2)$$

For a finite set  $\mathcal{D}$  of real numbers, with slight abuse of notation, we let  $Q_\alpha(\mathcal{D})$  be the  $\alpha$ -quantile of  $\mathcal{D}$ ; more precisely, we define  $Q_\alpha(\mathcal{D}) := Q_\alpha(X_{\mathcal{D}})$ , where  $X_{\mathcal{D}}$  denotes a random variable uniformly drawn from  $\mathcal{D}$ .

## 2 Algorithm

This section formalizes our algorithm design by specifying each stage of our meta-algorithm with a concrete procedure (namely, Algorithms 1, 2, 3 for Stages 1, 2, 3, respectively). It is worth noting that these are definitely not the only choices; in fact, an advantage of our meta-algorithm is its flexibility and modularity, in the sense that one can plug in different sub-routines to address various models and assumptions.

Before continuing, we introduce more notation that will be used throughout. For any  $1 \leq k \leq K$ , define

$$p_k := \frac{|\Omega_k^*|}{N} \quad \text{and} \quad r_k := \text{rank}(\mathbf{M}_k^*), \quad (3)$$

which represent the fraction of samples associated with the  $k$ -th component and the rank of the  $k$ -th ground-truth matrix  $\mathbf{M}_k^*$ , respectively. In addition, let the compact singular value decomposition (SVD) of  $\{\mathbf{M}_k^*\}$  be

$$\mathbf{M}_k^* = \mathbf{U}_k^* \boldsymbol{\Sigma}_k^* \mathbf{V}_k^{*\top}, \quad 1 \leq k \leq K, \quad (4)$$

where  $\mathbf{U}_k^* \in \mathbb{R}^{n_1 \times r_k}$  and  $\mathbf{V}_k^* \in \mathbb{R}^{n_2 \times r_k}$  consist of orthonormal columns, and  $\boldsymbol{\Sigma}_k^*$  is a diagonal matrix.

### 2.1 Stage 1: subspace estimation via a spectral method

**Procedure.** We propose to estimate the following joint column and row spaces:

$$\mathbf{U}^* := \text{col}\{\mathbf{U}_1^*, \dots, \mathbf{U}_K^*\} \quad \text{and} \quad \mathbf{V}^* := \text{col}\{\mathbf{V}_1^*, \dots, \mathbf{V}_K^*\} \quad (5)$$

by means of a spectral method. More specifically, we start by forming a data matrix

$$\mathbf{Y} := \frac{1}{N} \sum_{i=1}^N y_i \mathbf{A}_i, \quad (6)$$

and set  $\mathbf{U} \in \mathbb{R}^{n_1 \times R}$  (resp.  $\mathbf{V} \in \mathbb{R}^{n_2 \times R}$ ) to be a matrix whose columns consist of the top- $R$  left (resp. right) singular vectors of  $\mathbf{Y}$ , where

$$R := \text{rank}(\mathbb{E}[\mathbf{Y}]). \quad (7)$$

This method is summarized in Algorithm 1.

**Rationale.** To see why this might work, note that if  $\{\mathbf{A}_i\}$  consist of i.i.d. standard Gaussian entries, then

$$\begin{aligned} \mathbb{E}[\mathbf{Y}] &= \sum_{k=1}^K p_k \mathbb{E}[\langle \mathbf{A}_i, \mathbf{M}_k^* \rangle \mathbf{A}_i] = \sum_{k=1}^K p_k \mathbf{M}_k^* = \sum_{k=1}^K p_k \mathbf{U}_k^* \boldsymbol{\Sigma}_k^* \mathbf{V}_k^{*\top} \\ &= [\mathbf{U}_1^*, \mathbf{U}_2^*, \dots, \mathbf{U}_K^*] \begin{bmatrix} p_1 \boldsymbol{\Sigma}_1^* & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & p_2 \boldsymbol{\Sigma}_2^* & \mathbf{0} & \vdots \\ \vdots & \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \dots & \mathbf{0} & p_K \boldsymbol{\Sigma}_K^* \end{bmatrix} \begin{bmatrix} \mathbf{V}_1^{*\top} \\ \mathbf{V}_2^{*\top} \\ \vdots \\ \mathbf{V}_K^{*\top} \end{bmatrix}. \end{aligned} \quad (8)$$

Recalling the definitions of  $\mathbf{U}^*$  and  $\mathbf{V}^*$  in (5), we have

$$\text{col}\{\mathbb{E}[\mathbf{Y}]\} = \mathbf{U}^*, \quad \text{row}\{\mathbb{E}[\mathbf{Y}]\} = \mathbf{V}^*, \quad \text{rank}(\mathbf{U}^*) = \text{rank}(\mathbf{V}^*) = R$$

under some mild conditions (detailed in Section 3). This motivates the development of Algorithm 1.

### 2.2 Stage 2: initialization via low-dimensional mixed linear regression

**Key observations.** Suppose that there is an oracle informing us of the subspaces  $\mathbf{U}^*$  and  $\mathbf{V}^*$  defined in (5). Recognizing the basic relation  $\mathbf{M}_k^* = \mathbf{U}^* \mathbf{U}^{*\top} \mathbf{M}_k^* \mathbf{V}^* \mathbf{V}^{*\top}$  and defining

$$\mathbf{S}_k^* := \mathbf{U}^{*\top} \mathbf{M}_k^* \mathbf{V}^* \in \mathbb{R}^{R \times R}, \quad 1 \leq k \leq K, \quad (9)$$

---

**Algorithm 1:** Subspace estimation via a spectral method
 

---

- 1 **Input:** samples  $\{\mathbf{A}_i, y_i\}_{1 \leq i \leq N}$ , rank  $R$ .
  - 2 Compute  $\mathbf{Y} \leftarrow \frac{1}{N} \sum_{i=1}^N y_i \mathbf{A}_i$ .
  - 3 Let  $\mathbf{U} \in \mathbb{R}^{n_1 \times R}$  (resp.  $\mathbf{V} \in \mathbb{R}^{n_2 \times R}$ ) be the matrix consisting of the top- $R$  left (resp. right) singular vectors of  $\mathbf{Y}$ .
  - 4 **Output:**  $\mathbf{U}, \mathbf{V}$ .
- 

---

**Algorithm 2:** Initialization via low-dimensional mixed linear regression
 

---

- 1 **Input:** samples  $\{\mathbf{A}_i, y_i\}_{1 \leq i \leq N}$ , subspaces  $\mathbf{U}, \mathbf{V}$ , ranks  $\{r_k\}_{1 \leq k \leq K}$ .
  - 2 Transform  $\mathbf{a}_i \leftarrow \text{vec}(\mathbf{U}^\top \mathbf{A}_i \mathbf{V}), 1 \leq i \leq N$ .
  - 3 Obtain  $\{\widehat{\boldsymbol{\beta}}_k\}_{1 \leq k \leq K} \leftarrow$  the output of a black-box mixed linear regression solver (i.e. Algorithm 5) on  $\{\mathbf{a}_i, y_i\}_{1 \leq i \leq N}$ .
  - 4 **for**  $k = 1, \dots, K$  **do**
  - 5      $\mathbf{U}_k \boldsymbol{\Sigma}_k \mathbf{V}_k^\top \leftarrow$  rank- $r_k$  SVD of  $\mathbf{U} \widehat{\mathbf{S}}_k \mathbf{V}^\top$ , where  $\widehat{\mathbf{S}}_k := \text{mat}(\widehat{\boldsymbol{\beta}}_k)$ .
  - 6      $\mathbf{L}_k \leftarrow \mathbf{U}_k \boldsymbol{\Sigma}_k^{1/2}, \mathbf{R}_k \leftarrow \mathbf{V}_k \boldsymbol{\Sigma}_k^{1/2}$ .
  - 7 **Output:**  $\{\mathbf{L}_k, \mathbf{R}_k\}_{1 \leq k \leq K}$ .
- 

we can rewrite the measurements in hand as follows:

$$y_i = \begin{cases} \langle \mathbf{A}_i, \mathbf{M}_1^* \rangle = \langle \mathbf{A}_i, \mathbf{U}^* \mathbf{U}^{*\top} \mathbf{M}_1^* \mathbf{V}^* \mathbf{V}^{*\top} \rangle = \langle \mathbf{U}^{*\top} \mathbf{A}_i \mathbf{V}^*, \mathbf{S}_1^* \rangle, & \text{if } i \in \Omega_1^*, \\ \dots & \dots \\ \langle \mathbf{A}_i, \mathbf{M}_K^* \rangle = \langle \mathbf{U}^{*\top} \mathbf{A}_i \mathbf{V}^*, \mathbf{S}_K^* \rangle, & \text{if } i \in \Omega_K^*. \end{cases} \quad (10)$$

In other words, the presence of the oracle effectively reduces the original problem into a mixed linear regression problem in lower dimensions — that is, the problem of recovering  $\{\mathbf{S}_k^*\}$  from mixed linear measurements. If  $\{\mathbf{S}_k^*\}$  can be reliably estimated, then one can hope to recover  $\{\mathbf{M}_k^*\}$  via the following relation:

$$\mathbf{M}_k^* = \mathbf{U}^* \mathbf{U}^{*\top} \mathbf{M}_k^* \mathbf{V}^* \mathbf{V}^{*\top} = \mathbf{U}^* \mathbf{S}_k^* \mathbf{V}^{*\top}. \quad (11)$$

**Procedure.** While we certainly have no access to the aforementioned oracle in reality, Stage 1 described above provides us with subspace estimates  $\mathbf{U}$  and  $\mathbf{V}$  satisfying  $\mathbf{U}\mathbf{U}^\top \approx \mathbf{U}^* \mathbf{U}^{*\top}$  and  $\mathbf{V}\mathbf{V}^\top \approx \mathbf{V}^* \mathbf{V}^{*\top}$ . Treating these as surrogates of  $(\mathbf{U}^*, \mathbf{V}^*)$  (so that  $\mathbf{M}_k^* \approx \mathbf{U}\mathbf{U}^\top \mathbf{M}_k^* \mathbf{V}\mathbf{V}^\top$ ), we can view the measurements as

$$y_i = \begin{cases} \langle \mathbf{A}_i, \mathbf{M}_1^* \rangle \approx \langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top \mathbf{M}_1^* \mathbf{V}\mathbf{V}^\top \rangle = \langle \mathbf{U}^\top \mathbf{A}_i \mathbf{V}, \mathbf{S}_1 \rangle = \langle \mathbf{a}_i, \boldsymbol{\beta}_1 \rangle, & \text{if } i \in \Omega_1^*, \\ \dots & \dots \\ \langle \mathbf{A}_i, \mathbf{M}_K^* \rangle \approx \langle \mathbf{a}_i, \boldsymbol{\beta}_K \rangle, & \text{if } i \in \Omega_K^*, \end{cases} \quad (12)$$

which are mixed linear measurements about the following vectors/matrices:

$$\boldsymbol{\beta}_k := \text{vec}(\mathbf{S}_k) \in \mathbb{R}^{R^2}, \quad \mathbf{S}_k := \mathbf{U}^\top \mathbf{M}_k^* \mathbf{V} \in \mathbb{R}^{R \times R}, \quad 1 \leq k \leq K. \quad (13)$$

Here, the equivalent sensing vectors are defined to be  $\mathbf{a}_i := \text{vec}(\mathbf{U}^\top \mathbf{A}_i \mathbf{V}) \in \mathbb{R}^{R^2}$  for any  $1 \leq i \leq N$ . All this motivates us to resort to mixed linear regression algorithms for recovering  $\{\boldsymbol{\beta}_k\}$ . The proposed algorithm thus entails the following steps, with the precise procedure summarized in Algorithm 2.

- Invoke any mixed linear regression algorithm to obtain estimates  $\{\widehat{\boldsymbol{\beta}}_k\}_{1 \leq k \leq K}$  for  $\{\boldsymbol{\beta}_k\}_{1 \leq k \leq K}$  (up to global permutation). For concreteness, the current paper applies the tensor method (Algorithm 5) originally proposed in [YCS16]; this is a polynomial-time algorithm, with details deferred to Appendix A. To simplify presentation, let us assume here that the global permutation happens to be an identity map, so that  $\widehat{\boldsymbol{\beta}}_k$  is indeed a faithful estimate of  $\boldsymbol{\beta}_k$  ( $1 \leq k \leq K$ ). By simple matricization,  $\widehat{\boldsymbol{\beta}}_k$  leads to a reliable estimate  $\widehat{\mathbf{S}}_k$  of  $\mathbf{S}_k$ .

---

**Algorithm 3:** Scaled Truncated Gradient Descent (ScaledTGD) for recovering  $\mathbf{M}_k^*$ 


---

1 **Input:** samples  $\{\mathbf{A}_i, y_i\}_{1 \leq i \leq N}$ , initialization  $\mathbf{L}^0 \in \mathbb{R}^{n_1 \times r_k}$ ,  $\mathbf{R}^0 \in \mathbb{R}^{n_2 \times r_k}$ , step size  $\eta$ , truncating fraction  $\alpha$ .

2 **for**  $t = 0, 1, 2, \dots, T_0 - 1$  **do**

3

$$\mathbf{L}^{t+1} \leftarrow \mathbf{L}^t - \frac{\eta}{N} \sum_{i \in \Omega^t} (\langle \mathbf{A}_i, \mathbf{L}^t(\mathbf{R}^t)^\top \rangle - y_i) \mathbf{A}_i \mathbf{R}^t ((\mathbf{R}^t)^\top \mathbf{R}^t)^{-1},$$

$$\mathbf{R}^{t+1} \leftarrow \mathbf{R}^t - \frac{\eta}{N} \sum_{i \in \Omega^t} (\langle \mathbf{A}_i, \mathbf{L}^t(\mathbf{R}^t)^\top \rangle - y_i) \mathbf{A}_i^\top \mathbf{L}^t ((\mathbf{L}^t)^\top \mathbf{L}^t)^{-1},$$

where  $\Omega^t := \{1 \leq i \leq N : |\langle \mathbf{A}_i, \mathbf{L}^t(\mathbf{R}^t)^\top \rangle - y_i| \leq \tau_t\}$ ,  $\tau_t := Q_\alpha(\{|\langle \mathbf{A}_i, \mathbf{L}^t(\mathbf{R}^t)^\top \rangle - y_i|\}_{1 \leq i \leq N})$ .

4 **Output:**  $\mathbf{L}^{T_0}, \mathbf{R}^{T_0}$ .

---

- Given the observation that

$$\mathbf{U} \widehat{\Sigma}_k \mathbf{V}^\top \approx \mathbf{U} \Sigma_k \mathbf{V}^\top = \mathbf{U} \mathbf{U}^\top \mathbf{M}_k^* \mathbf{V} \mathbf{V}^\top \approx \mathbf{U}^* \mathbf{U}^{*\top} \mathbf{M}_k^* \mathbf{V}^* \mathbf{V}^{*\top} = \mathbf{M}_k^*, \quad (14)$$

we propose to compute the rank- $r_k$  SVD — denoted by  $\mathbf{U}_k \Sigma_k \mathbf{V}_k^\top$  — of the matrix  $\mathbf{U} \widehat{\Sigma}_k \mathbf{V}^\top$  for each  $1 \leq k \leq K$ . This in turn leads to our initial estimate for the low-rank factors

$$\mathbf{L}_k := \mathbf{U}_k \Sigma_k^{1/2} \in \mathbb{R}^{n_1 \times r_k}, \quad \text{and} \quad \mathbf{R}_k := \mathbf{V}_k \Sigma_k^{1/2} \in \mathbb{R}^{n_2 \times r_k}. \quad (15)$$

### 2.3 Stage 3: local refinement via scaled truncated gradient descent (ScaledTGD)

Suppose that an initial point  $\mathbf{L}^0(\mathbf{R}^0)^\top$  lies within a reasonably small neighborhood of  $\mathbf{M}_k^*$  for some  $1 \leq k \leq K$ . Stage 3 serves to locally refine this initial estimate, moving it closer to our target  $\mathbf{M}_k^*$ . Towards this end, we propose to deploy the following update rule termed *scaled truncated gradient descent* (ScaledTGD):

$$\mathbf{L}^{t+1} = \mathbf{L}^t - \frac{\eta}{N} \sum_{i \in \Omega^t} (\langle \mathbf{A}_i, \mathbf{L}^t(\mathbf{R}^t)^\top \rangle - y_i) \mathbf{A}_i \mathbf{R}^t ((\mathbf{R}^t)^\top \mathbf{R}^t)^{-1}, \quad (16a)$$

$$\mathbf{R}^{t+1} = \mathbf{R}^t - \frac{\eta}{N} \sum_{i \in \Omega^t} (\langle \mathbf{A}_i, \mathbf{L}^t(\mathbf{R}^t)^\top \rangle - y_i) \mathbf{A}_i^\top \mathbf{L}^t ((\mathbf{L}^t)^\top \mathbf{L}^t)^{-1}, \quad (16b)$$

where  $\eta > 0$  denotes the step size. Here,  $\Omega^t \subseteq \{1, 2, \dots, N\}$  is an adaptive and iteration-varying index set designed to mimic the index set  $\Omega_k^*$ . Indeed, if  $\Omega^t = \Omega_k^*$ , the aforementioned update rule reduces to the ScaledGD method developed for vanilla low-rank matrix sensing (see [TMC20]), which is guaranteed to converge to  $\mathbf{M}_k^*$  in the presence of a suitable initialization. Here, the rescaling matrix  $((\mathbf{R}^t)^\top \mathbf{R}^t)^{-1}$  (resp.  $((\mathbf{L}^t)^\top \mathbf{L}^t)^{-1}$ ) acts as a pre-conditioner of the conventional gradient  $\sum_{i \in \Omega^t} (\langle \mathbf{A}_i, \mathbf{L}^t(\mathbf{R}^t)^\top \rangle - y_i) \mathbf{A}_i \mathbf{R}^t$  (resp.  $\sum_{i \in \Omega^t} (\langle \mathbf{A}_i, \mathbf{L}^t(\mathbf{R}^t)^\top \rangle - y_i) \mathbf{A}_i^\top \mathbf{L}^t$ ), which effectively accelerates convergence when  $\mathbf{M}_k^*$  is ill-conditioned. See [TMC20] for more intuitions and justifications of this rescaling strategy.

Viewed in this light, the key to ensuring effectiveness of ScaledTGD lies in the design of the index set  $\Omega^t$ . If we know *a priori* that  $\mathbf{L}^t(\mathbf{R}^t)^\top \approx \mathbf{M}_k^*$ , then it is intuitively clear that  $|\langle \mathbf{A}_i, \mathbf{L}^t(\mathbf{R}^t)^\top \rangle - y_i|$  typically has a smaller scale for a sample  $i \in \Omega_k^*$  when compared with those  $i \notin \Omega_k^*$ . This motivates us to include in  $\Omega^t$  a certain fraction (denoted by  $0 < \alpha < 1$ ) of samples enjoying the smallest empirical loss  $|\langle \mathbf{A}_i, \mathbf{L}^t(\mathbf{R}^t)^\top \rangle - y_i|$ . Intuitively, the fraction  $\alpha$  should not be too large in which case  $\Omega^t$  is likely to contain samples outside  $\Omega_k^*$ ; on the other hand,  $\alpha$  should not be chosen too small in order not to waste information. As it turns out, choosing  $0.6p_k \leq \alpha \leq 0.8p_k$  strikes a suitable balance and works well for our purpose. See Algorithm 3 for a precise description.

### 2.4 The full algorithm

With the three stages fully described, we can specify the whole algorithm in Algorithm 4, with the choices of algorithmic parameters listed in Table 1. Note that the discussion in Section 2.3 focuses on estimating a

---

**Algorithm 4:** A fully specified three-stage algorithm for mixed matrix sensing
 

---

- 1 **Input:** independent samples  $\{\mathbf{A}_i, y_i\}_{1 \leq i \leq N}$ ,  $\{\mathbf{A}'_i, y'_i\}_{1 \leq i \leq N_{\text{MLR}}}$ , parameters  $R, \{r_k, \eta_k, \alpha_k\}_{1 \leq k \leq K}$  (see Table 1).
  - 2 Run Algorithm 1 with  $\{\mathbf{A}_i, y_i\}_{1 \leq i \leq N}$  and  $R$  to obtain  $\mathbf{U}, \mathbf{V}$ .
  - 3 Run Algorithm 2 with  $\{\mathbf{A}'_i, y'_i\}_{1 \leq i \leq N_{\text{MLR}}}$ ,  $\mathbf{U}, \mathbf{V}$  and  $\{r_k\}_{1 \leq k \leq K}$  to obtain  $\{\mathbf{L}_k, \mathbf{R}_k\}_{1 \leq k \leq K}$ .
  - 4 **for**  $k = 1, 2, \dots, K$  **do**
  - 5     Run Algorithm 3 on  $\{\mathbf{A}_i, y_i\}_{1 \leq i \leq N}$  with  $(\mathbf{L}^0, \mathbf{R}^0) \leftarrow (\mathbf{L}_k, \mathbf{R}_k), \eta_k, \alpha_k$  to obtain  $\mathbf{L}^{T_0}, \mathbf{R}^{T_0}$ .
  - 6     Set  $\mathbf{M}_k \leftarrow \mathbf{L}^{T_0} (\mathbf{R}^{T_0})^\top$ .
  - 7 **Output:**  $\{\mathbf{M}_k\}_{1 \leq k \leq K}$ .
- 

Table 1: Our choices of the algorithmic parameters in Algorithm 4.

Algorithm 1	Rank $R = \text{rank}(\sum_k p_k \mathbf{M}_k^*)$ .
Algorithm 2	Ranks $r_k = \text{rank}(\mathbf{M}_k^*), 1 \leq k \leq K$ .
Algorithm 3 (for $\mathbf{M}_k^*$ )	Step size $0 < \eta_k \leq 1.3/p_k$ , truncating fraction $0.6p_k \leq \alpha_k \leq 0.8p_k$ .

single component; in order to recover all  $K$  components  $\{\mathbf{M}_k^*\}_{1 \leq k \leq K}$ , we simply need to run Algorithm 3 for  $K$  times (which can be executed in parallel). In addition, Algorithm 4 is built upon sample splitting: while Stages 1 and 3 employ the same set of samples  $\{\mathbf{A}_i, y_i\}_{1 \leq i \leq N}$ , Stage 2 (i.e. Line 3 of Algorithm 4) operates upon an *independent* set of samples  $\{\mathbf{A}'_i, y'_i\}_{1 \leq i \leq N_{\text{MLR}}}$  (where ‘‘MLR’’ stands for ‘‘mixed linear regression’’), thus resulting in a total sample complexity of  $N + N_{\text{MLR}}$ . The main purpose of sample splitting is to decouple statistical dependency across stages and facilitate analysis. Finally, the interested reader is referred to Appendix D for a discussion regarding how to estimate certain parameters in Algorithm 4 if they are not known *a priori*.

### 3 Main results

#### 3.1 Models and assumptions

For notational convenience, let us define the following parameters:

$$n := \max\{n_1, n_2\}, \quad r := \max_{1 \leq k \leq K} r_k, \quad \kappa := \max_{1 \leq k \leq K} \kappa(\mathbf{M}_k^*), \quad \text{and} \quad \Gamma := \frac{\max_{1 \leq k \leq K} \|\mathbf{M}_k^*\|_{\text{F}}}{\min_{1 \leq k \leq K} \|\mathbf{M}_k^*\|_{\text{F}}}, \quad (17)$$

where  $\kappa(\mathbf{M}_k^*) := \sigma_1(\mathbf{M}_k^*)/\sigma_{r_k}(\mathbf{M}_k^*)$  stands for the condition number of  $\mathbf{M}_k^*$ . This paper focuses on the *Gaussian design*, where the entries of each design matrix  $\mathbf{A}_i$  are independently drawn from the standard Gaussian distribution. In addition, we assume that the samples drawn from the  $K$  components are reasonably *well-balanced* in the sense that for all  $1 \leq k \leq K$ ,

$$p_k = \frac{|\Omega_k^*|}{N} \asymp \frac{1}{K}, \quad (18)$$

where  $\Omega_k^*$  is the index set for the  $k$ -th component (see (1)). We assume that this well-balancedness assumption holds for both sets of samples  $\{\mathbf{A}_i, y_i\}_{1 \leq i \leq N}$  and  $\{\mathbf{A}'_i, y'_i\}_{1 \leq i \leq N_{\text{MLR}}}$ .

Next, we introduce an incoherence parameter that plays a crucial role in our theoretical development.

**Definition 1.** *The incoherence parameter  $\mu \geq 0$  is the smallest quantity that satisfies*

$$\|\mathbf{U}_i^{*\top} \mathbf{U}_j^*\|_{\text{F}} \leq \frac{\mu r}{\sqrt{n_1}}, \quad \text{and} \quad \|\mathbf{V}_i^{*\top} \mathbf{V}_j^*\|_{\text{F}} \leq \frac{\mu r}{\sqrt{n_2}} \quad \text{for all } 1 \leq i < j \leq K. \quad (19)$$

The incoherence parameter  $\mu$  takes value on  $[0, \sqrt{n/r}]$ . As an example, if  $\{\mathbf{U}_k^*\}_{1 \leq k \leq K}$  (resp.  $\{\mathbf{V}_k^*\}_{1 \leq k \leq K}$ ) are random low-dimensional subspaces in  $\mathbb{R}^{n_1}$  (resp.  $\mathbb{R}^{n_2}$ ), then for any  $i \neq j$ ,  $\|\mathbf{U}_i^{*\top} \mathbf{U}_j^*\|_{\text{F}}$  (resp.  $\|\mathbf{V}_i^{*\top} \mathbf{V}_j^*\|_{\text{F}}$ )

is on the order of  $\sqrt{r_i r_j / n_1}$  (resp.  $\sqrt{r_i r_j / n_2}$ ), which is further upper bounded by  $r / \sqrt{n_1}$  (resp.  $r / \sqrt{n_2}$ ). This observation motivates our definition of the incoherence parameter. One of our main technical assumptions is that the column (resp. row) spaces of the ground-truth matrices are mutually weakly correlated — defined through the parameter  $\mu$  — which covers a broad range of settings.

**Assumption 1.** *The incoherence parameter  $\mu$  is upper bounded by*

$$\mu \leq \frac{\sqrt{\min\{n_1, n_2\}}}{2r \max\{K, \sqrt{K}\Gamma}}. \quad (20)$$

## 3.2 Theoretical guarantees

**Exact recovery in the absence of noise.** Our first main result uncovers that, in the noiseless case, Algorithm 4 achieves exact recovery efficiently, in terms of both sample and computational complexities.

**Theorem 1** (Exact recovery). *Consider the noiseless case (1) under the assumptions in Section 3.1. Suppose*

$$N \geq C_1 K^3 r^2 \kappa^2 \Gamma^2 \max\{K^2 \Gamma^4, r \kappa^2\} \cdot n \log N \quad \text{and} \quad N_{\text{MLR}} \geq C_2 K^8 r^2 \Gamma^{12} \max\{K^2, r \kappa^2\} \cdot \log n \cdot \log^3 N_{\text{MLR}} \quad (21)$$

for some sufficiently large constants  $C_1, C_2 > 0$ . Then with probability at least  $1 - o(1)$ , there exists some permutation  $\pi : \{1, \dots, K\} \mapsto \{1, \dots, K\}$  such that the outputs of Algorithm 4 obey for all  $1 \leq k \leq K$

$$\|\mathbf{M}_{\pi(k)} - \mathbf{M}_k^*\|_{\text{F}} \leq (1 - c_0 \eta_k p_k)^{T_0} \|\mathbf{M}_k^*\|_{\text{F}} \quad (22)$$

for some universal constant  $0 < c_0 < 1/4$ , where  $T_0$  is the number of iterations used in Algorithm 3.

The proof can be found in Section 5. Two implications are in order.

- Suppose that the parameters  $K, r, \kappa, \Gamma = O(1)$ . In order to achieve exact recovery, the sample size  $N$  in (21) only needs to scale as  $O(n \log n)$ , while  $N_{\text{MLR}}$  only needs to exceed the order of  $\log^4 n$ .
- By setting the step size  $\eta_k = c_1 / p_k$  for some constant  $0 < c_1 \leq 1.3$ , we see that the third stage (i.e. ScaledTGD) achieves linear convergence with a *constant* contraction rate, which is independent of the condition number  $\kappa(\mathbf{M}_k^*)$  of the matrix  $\mathbf{M}_k^*$ .

**Stability vis-à-vis noise.** Moving on to the more realistic case with noise, we consider the following set of samples  $\{\mathbf{A}_i, y_i\}_{1 \leq i \leq N}$ :

$$\zeta_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2), \quad y_i = \begin{cases} \langle \mathbf{A}_i, \mathbf{M}_1^* \rangle + \zeta_i, & \text{if } i \in \Omega_1^*, \\ \dots & \\ \langle \mathbf{A}_i, \mathbf{M}_K^* \rangle + \zeta_i, & \text{if } i \in \Omega_K^*. \end{cases} \quad (23)$$

The set  $\{\mathbf{A}'_i, y'_i\}_{1 \leq i \leq N_{\text{MLR}}}$  is independently generated in a similar manner. Our next result reveals that the proposed algorithm is stable against Gaussian noise. The proof is postponed to Section 5.

**Theorem 2** (Stable recovery). *Consider the noisy model (23) under the assumptions of Section 3.1. Suppose that the sample sizes satisfy (21), and that the noise level satisfies*

$$\sigma \leq c \min_{1 \leq k \leq K} \|\mathbf{M}_k^*\|_{\text{F}} \cdot \min \left\{ \frac{1}{K}, \frac{1}{\sqrt{r\kappa}} \right\} \quad (24)$$

for some sufficiently small constant  $c > 0$ . Then with probability at least  $1 - o(1)$ , there exists some permutation  $\pi : \{1, \dots, K\} \mapsto \{1, \dots, K\}$  such that the outputs of Algorithm 4 obey for all  $1 \leq k \leq K$

$$\|\mathbf{M}_{\pi(k)} - \mathbf{M}_k^*\|_{\text{F}} \leq (1 - c_0 \eta_k p_k)^{T_0} \|\mathbf{M}_k^*\|_{\text{F}} + C_0 \max \left\{ \sigma \sqrt{\frac{nrK^3 \log N}{N}}, \frac{K\sigma^2}{\min_{j:j \neq k} \|\mathbf{M}_j^* - \mathbf{M}_k^*\|_{\text{F}}} \right\}, \quad (25)$$

where  $0 < c_0 < 1/4$  and  $C_0 > 0$  are some universal constants, and  $T_0$  is the number of iterations used in Algorithm 3.



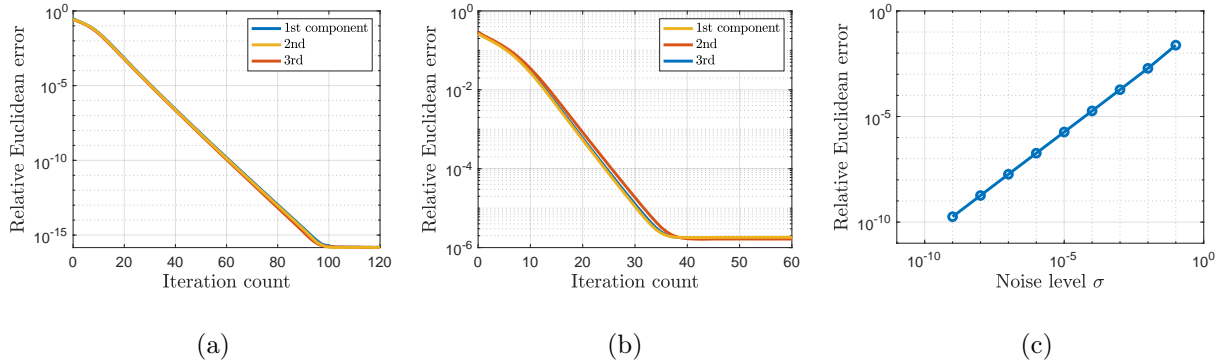


Figure 1: (a) The relative Euclidean error vs. the iteration count of ScaledTGD in Stage 3 of Algorithm 4 for each of the three components, in the noiseless case. (b) Convergence of ScaledTGD in the noisy case  $\sigma = 10^{-5}$ . (c) The largest relative Euclidean error (at convergence) of ScaledTGD in Algorithm 4, vs. the noise level  $\sigma$ . Each data point is an average over 10 independent trials.

Theorem 2 asserts that, when initialized using the proposed schemes, the ScaledTGD algorithm converges linearly until an error floor is hit. To interpret the statistical guarantees (25), we find it helpful to define the signal-to-noise-ratio (SNR) w.r.t.  $\mathbf{M}_k^*$  as follows:

$$\text{SNR}_k := \frac{\mathbb{E} \left[ |\langle \mathbf{A}_i, \mathbf{M}_k^* \rangle|^2 \right]}{\mathbb{E}[\zeta_i^2]} = \frac{\|\mathbf{M}_k^*\|_F^2}{\sigma^2}. \quad (26)$$

This together with the simple consequence  $\min_{j:j \neq k} \|\mathbf{M}_j^* - \mathbf{M}_k^*\|_F \gtrsim \|\mathbf{M}_k^*\|_F$  of Assumption 1 implies that

$$\frac{\|\mathbf{M}_{\pi(k)} - \mathbf{M}_k^*\|_F}{\|\mathbf{M}_k^*\|_F} \lesssim \max \left\{ \frac{1}{\sqrt{\text{SNR}_k}} \sqrt{\frac{nrK^3 \log N}{N}}, \frac{K}{\text{SNR}_k} \right\} \quad (27)$$

as long as the iteration number  $T_0$  is sufficiently large. Here, the first term on the right-hand side of (27) matches the *minimax lower bound* for low-rank matrix sensing [CP11, Theorem 2.5] (the case with  $K = 1$ ) up to a factor of  $K\sqrt{\log N}$ . In contrast, the second term on the right-hand side of (27) — which becomes very small as  $\text{SNR}_k$  grows — is not a function of the sample size  $N$  and does not vanish as  $N \rightarrow \infty$ . This term arises since, even at the population level, the point  $(\mathbf{L}, \mathbf{R})$  satisfying  $\mathbf{L}\mathbf{R}^\top = \mathbf{M}_k^*$  is not a fixed point of the ScaledTGD update rule, due to the presence of mislabeled samples.

### 3.3 Numerical experiments

To validate our theoretical findings, we conduct a series of numerical experiments. To match practice, we do not deploy sample splitting (given that it is merely introduced to simplify analysis), and reuse the same dataset of size  $N$  for all three stages. Throughout the experiments, we set  $n_1 = n_2 = n = 120$ ,  $r = 2$ , and  $K = 3$ . For each  $k$ , we let  $p_k = 1/K$  and  $\boldsymbol{\Sigma}_k^* = \mathbf{I}_r$ , and generate  $\mathbf{U}_k^*$  and  $\mathbf{V}_k^*$  as random  $r$ -dimensional subspaces in  $\mathbb{R}^n$ . We fix the sample size to be  $N = 90nrK$ . The algorithmic parameters are chosen according to our recommendations in Table 1. For instance, for each run of ScaledTGD, we set the step size as  $\eta = 1.3K$  and the truncation fraction as  $\alpha = 0.8/K$ .

**Linear convergence of ScaledTGD.** Our first series of experiments aims at verifying the linear convergence of ScaledTGD towards the ground-truth matrices  $\{\mathbf{M}_k^*\}$  when initialized using the outputs of Stage 2. We consider both the noiseless case (i.e.  $\sigma = 0$ ) and the noisy case  $\sigma = 10^{-5}$ . Figures 1(a) and 1(b) plot the relative Euclidean error  $\|\mathbf{L}^t(\mathbf{R}^t)^\top - \mathbf{M}_k^*\|_F / \|\mathbf{M}_k^*\|_F$  versus the iteration count  $t$  for each component  $1 \leq k \leq 3$ . It is easily seen from Figures 1(a) and 1(b) that ScaledTGD, when seeded with the outputs from Stage 2, converges linearly to the ground-truth matrices  $\{\mathbf{M}_k^*\}$  in the absence of noise, and to within a small neighborhood of  $\{\mathbf{M}_k^*\}$  in the noisy setting.

**Estimation error in the presence of random noise.** The second series of experiments investigates the stability of the three-stage algorithm in the presence of random noise. We vary the noise level within  $[10^{-9}, 10^{-1}]$ . Figure 1(c) plots the largest relative Euclidean error  $\max_{1 \leq k \leq K} \|\mathbf{M}_k - \mathbf{M}_k^*\|_{\text{F}} / \|\mathbf{M}_k^*\|_{\text{F}}$  (where  $\{\mathbf{M}_k\}$  are the outputs of Algorithm 4) versus the noise level  $\sigma$ , showing that the recovering error is indeed linear in  $\sigma$ , as predicted by our theory.

## 4 Prior work

**Low-rank matrix recovery.** There exists a vast literature on low-rank matrix recovery (e.g. [CR09, KMO10, BNS16, CC17, MWCC20, CCFM19, SL16, CLS15, JNS13, CCF<sup>+</sup>ar, CFMY19, SQW18, CLL20, DC20, NNS<sup>+</sup>14, CCG15, CCD<sup>+</sup>19, ACHL19, ZQW20, ZWYG18, LMZ18, PKCS17]); we refer the readers to [CC18b, CLC19] for an overview of this extensively studied topic. Most related to our work is the problem of matrix sensing (or low-rank matrix recovery from linear measurements). While convex relaxation [CR09, RFP10, CP11] enjoys optimal statistical performance, two-stage non-convex approaches [ZL15, TBS<sup>+</sup>16, TMC20] have received growing attention in recent years, due to their ability to achieve statistical and computational efficiency at once. Our three-stage algorithm is partially inspired by the two-stage approach along this line. It is worth mentioning that the non-convex loss function associated with low-rank matrix sensing enjoys benign landscape, which in turn enables tractable global convergence of simple first-order methods [BNS16, GJZ17, ZLTW18, LMZ18, LZT19].

**Mixed linear regression.** Being a classical problem in statistics [QR78], mixed linear regression has attracted much attention due to its broad applications in music perception [DV89, VT02], health care [DH00], trajectory clustering [GS99], plant science [Tur00], neuroscience [YPCR18], to name a few. While computationally intractable in the worst case [YCS14], mixed linear regression can be solved efficiently under certain statistical models on the design matrix. Take the two-component case for instance: efficient methods include alternating minimization with initialization via grid search [YCS14], EM with random initialization [KYB19, KQC<sup>+</sup>19], and convex reformulations [CYC14, HJ18], where EM further achieves minimax estimation guarantees [CYC14] in the presence of Gaussian noise [KHC20]. Mixed linear regression becomes substantially more challenging when the number  $K$  of components is allowed to grow with  $n$ . The state-of-the-art method — namely, the Fourier moment method [CLS20] — achieves sub-exponential sample and computational complexities w.r.t.  $K$ , whereas other methods (e.g. the method of moments [LL18] and grid search over  $K$ -dimensional subspaces [SS19a]) all have exponential dependence on  $K$ . It turns out that by restricting the ground-truth vectors to be in “general position” (e.g. linearly independent), tensor methods [YCS16, CL13, SJA16, ZJD16] solve mixed linear regression with polynomial sample and computational complexities in  $K$ . It is worth noting that most of the prior work focused on the Gaussian design for theoretical analysis, with a few exceptions [CYC14, HJ18, SS19a]. Another line of work [KC07, SBVDG10, YPCR18, KMMP19, MP20] considered mixed linear regression with sparsity, which is beyond the scope of the current paper.

**Mixed low-rank matrix estimation.** Moving beyond mixed linear regressions, there are a few papers that tackle mixtures of low-rank models. For example, [YC15] proposed a regularized EM algorithm and applied it to mixed matrix sensing with two symmetric components; however, only local convergence was investigated therein. Additionally, [PA18] was the first to systematically study mixed matrix completion, investigating the identifiability conditions and sample complexities of this problem; however, the heuristic algorithm proposed therein comes without provable guarantees.

**Iterative truncated loss minimization.** Least trimmed square [Rou84] is a classical method for robust linear regression. Combining the idea of trimming (i.e. selecting a subset of “good” samples) with iterative optimization algorithms (e.g. gradient descent and its variants) leads to a general paradigm of iterative truncated loss minimization — a principled method for improving robustness w.r.t. heavy-tailed data, adversarial outliers, etc. [SS19b, SWS20]. Successful applications of this kind include linear regression [BJK15], mixed linear regression [SS19a], phase retrieval [CC17, ZCL18], matrix sensing [LCZL20], and learning entangled single-sample distributions [YL20], among others.

**Multi-task learning and meta-learning.** The aim of multi-task learning [Car97, Bax00, BDS03, AZ05, EMP05, AEP07, JSRR10, PLW+20, PL14, MPRP16] is to simultaneously learn a model that connect multiple *related* tasks. Exploiting the similarity across tasks enables improved performance for learning each individual task, and leads to enhance generalization capabilities for unseen but related tasks with limited samples. This paradigm (or its variants) is also referred to in the literature as meta-learning [FAL17, TJJ20] (i.e. learning-to-learn), transfer learning [PY09], and few-shot learning [SSZ17, DHK+20], depending on the specific scenarios of interest. Our study on learning mixture of models is related to the probabilistic approach taken in multi-task learning and meta-learning, in which all the tasks (both the training and the testing ones) are independently sampled from a common environment, i.e. a prior distribution of tasks [Bax00]. See [KSS+20, KSKO20] for recent efforts that make explicit the connection between mixed linear regression and meta-learning.

## 5 Analysis

In this section, we present the proofs of Theorems 1 and 2. Our analysis is modular in the sense that we deliver the performance guarantees for the three stages separately that are independent of each other. For instance, one can replace the tensor method in Stage 2 by any other mixed linear regression solver with provable guarantees, without affecting Stages 1 and 3.

**Stage 1.** The first result confirms that given enough samples, Algorithm 1 outputs reasonable estimates of the subspaces  $(\mathbf{U}^*, \mathbf{V}^*)$  (cf. (5)). The proof is deferred to Appendix B.1.

**Theorem 3.** Consider the model (23) under the assumptions in Section 3.1. Recall the definitions of  $\kappa$  and  $\Gamma$  in (17). For any  $0 < \delta < 1$ , the estimates  $\mathbf{U}$  and  $\mathbf{V}$  returned by Algorithm 1 satisfy

$$\max \left\{ \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top}\|, \|\mathbf{V}\mathbf{V}^\top - \mathbf{V}^*\mathbf{V}^{*\top}\| \right\} \lesssim \delta K \sqrt{r} \kappa \left( \Gamma + \frac{1}{\sqrt{Kr}} \frac{\sigma}{\min_k \|\mathbf{M}_k^*\|_{\text{F}}} \right) \quad (28)$$

with probability at least  $1 - Ce^{-cn}$  for some universal constants  $C, c > 0$ , provided that the sample size obeys

$$N \geq C_0 \frac{nrK}{\delta^2} \log \frac{1}{\delta} \quad (29)$$

for some sufficiently large constant  $C_0 > 0$ .

**Stage 2.** Next, we demonstrate that the tensor method employed in Algorithm 2 reliably solves the intermediate mixed linear regression problem defined in (12). The proof is postponed to Appendix B.2.

**Theorem 4.** Consider the model (23) under the assumptions in Section 3.1. Suppose that the subspace estimates  $\mathbf{U}$  and  $\mathbf{V}$  are independent of  $\{\mathbf{A}_i, y_i\}_{1 \leq i \leq N}$  and obey  $\max\{\|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top}\|, \|\mathbf{V}\mathbf{V}^\top - \mathbf{V}^*\mathbf{V}^{*\top}\|\} \leq c_1/(K\Gamma^2)$  for some sufficiently small constant  $c_1 > 0$ . Let  $\{\hat{\beta}_k\}_{1 \leq k \leq K}$  be the estimates returned by Line 3 of Algorithm 2. Given any  $0 < \epsilon \leq c_2/K$ , there exists a permutation  $\pi(\cdot) : \{1, \dots, K\} \mapsto \{1, \dots, K\}$  such that

$$\|\hat{\beta}_{\pi(k)} - \beta_k\|_2 \leq \epsilon \cdot \max_{1 \leq j \leq K} \|\mathbf{M}_j^*\|_{\text{F}} \quad \text{for all } 1 \leq k \leq K \quad (30)$$

with probability at least  $1 - O(1/\log n)$ , provided that the sample size obeys

$$N \geq C \frac{K^8 r^2}{\epsilon^2} \left( \Gamma^{10} + \frac{\sigma^{10}}{\min_k \|\mathbf{M}_k^*\|_{\text{F}}^{10}} \right) \log n \cdot \log^3 N. \quad (31)$$

Here,  $c_2 > 0$  (resp.  $C > 0$ ) is some sufficiently small (resp. large) constant.

From now on, we shall assume without loss of generality that  $\pi(\cdot)$  is an identity map (i.e.  $\pi(k) = k$ ) to simplify the presentation. Our next result transfers the estimation error bounds for  $\mathbf{U}, \mathbf{V}$  and  $\{\hat{\beta}_k\}$  to that for  $\{\mathbf{L}_k \mathbf{R}_k^\top\}$ , thus concluding the analysis of Stage 2; see Appendix B.3 for a proof.

**Proposition 1.** The estimates  $\{\mathbf{L}_k, \mathbf{R}_k\}_{k=1}^K$  computed in Lines 4-6 of Algorithm 2 obey

$$\|\mathbf{L}_k \mathbf{R}_k^\top - \mathbf{M}_k^*\|_{\text{F}} \leq 2 \max \left\{ \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top}\|, \|\mathbf{V}\mathbf{V}^\top - \mathbf{V}^*\mathbf{V}^{*\top}\| \right\} \|\mathbf{M}_k^*\|_{\text{F}} + 2 \|\hat{\beta}_k - \beta_k\|_2 \quad (32)$$

for all  $1 \leq k \leq K$ .

**Stage 3.** The last result guarantees that Algorithm 3 — when suitably initialized — converges linearly towards  $\mathbf{M}_k^*$  up to a certain error floor. Here  $\mathbf{M}_k^*$  is the closest among  $\{\mathbf{M}_j^*\}_{1 \leq j \leq K}$  to the point  $\mathbf{L}^0(\mathbf{R}^0)^\top$ . The proof can be found in Appendix B.4.

**Theorem 5.** Consider the model (23) under the assumptions in Section 3.1. Suppose that the noise level obeys (24). Choose the step size  $\eta$  and truncating fraction  $\alpha$  such that  $0 < \eta \leq 1.3/p_k$  and  $0.6p_k \leq \alpha \leq 0.8p_k$ . Given any  $0 < \delta < c_0/K$ , if  $\mathbf{L}^0 \in \mathbb{R}^{n_1 \times r_k}$  and  $\mathbf{R}^0 \in \mathbb{R}^{n_2 \times r_k}$  obey

$$\|\mathbf{L}^0(\mathbf{R}^0)^\top - \mathbf{M}_k^*\|_{\text{F}} \leq c_1 \|\mathbf{M}_k^*\|_{\text{F}} \cdot \min \left\{ \frac{1}{\sqrt{r_k}}, \frac{1}{K} \right\}, \quad (33)$$

then with probability at least  $1 - Ce^{-cn}$  the iterates of Algorithm 3 satisfy

$$\|\mathbf{L}^t(\mathbf{R}^t)^\top - \mathbf{M}_k^*\|_{\text{F}} \leq (1 - c_2\eta p_k)^t \|\mathbf{L}^0(\mathbf{R}^0)^\top - \mathbf{M}_k^*\|_{\text{F}} + C_2 \max \left\{ K\sigma\delta, \frac{K\sigma^2}{\min_{j:j \neq k} \|\mathbf{M}_j^* - \mathbf{M}_k^*\|_{\text{F}}} \right\} \quad (34)$$

for all  $t \geq 0$ , provided that the sample size exceeds  $N \geq C_0 \frac{nr_k}{\delta^2} \log N$ . Here,  $0 < c_2 < 1/4$  and  $C, c, C_2 > 0$  are some universal constants, and  $c_0, c_1 > 0$  (resp.  $C_0 > 0$ ) are some sufficiently small (resp. large) constants.

**Putting pieces together: proof of Theorems 1 and 2.** With the above performance guarantees in place, we are ready to establish the main theorems. Note that due to sample splitting in Algorithm 4, we shall apply Theorems 3 and 5 to the dataset  $\{\mathbf{A}_i, y_i\}_{1 \leq i \leq N}$ , and Theorem 4 to the dataset  $\{\mathbf{A}'_i, y'_i\}_{1 \leq i \leq N_{\text{MLR}}}$ . Set

$$\delta \leq c_3 \frac{1}{K\sqrt{r_k}\Gamma} \min \left\{ \frac{1}{\sqrt{r_k}}, \frac{1}{K\Gamma^2} \right\}, \quad \text{and} \quad \epsilon \leq c_4 \frac{1}{\Gamma} \min \left\{ \frac{1}{\sqrt{r_k}}, \frac{1}{K} \right\},$$

for some sufficiently small constants  $c_3, c_4 > 0$  in Theorems 3 and 4. These choices — in conjunction with our assumption on  $\sigma$  in Theorem 2, as well as Proposition 1 — guarantee that the initialization  $\mathbf{L}^0(\mathbf{R}^0)^\top$  lies in the neighborhood of  $\mathbf{M}_k^*$  as required by (33). This allows us to invoke Theorem 5 to conclude the proof of Theorem 2. Finally, Theorem 1 follows by simply setting the noise level  $\sigma = 0$  in Theorem 2.

## 6 Discussion

This paper develops a three-stage algorithm for the mixed low-rank matrix sensing problem, which is provably efficient in terms of both sample and computational complexities. Having said this, there are numerous directions that are worthy of further investigations; we single out a few in the following.

To begin with, while our required sample complexity scales linearly (and optimally) w.r.t. the matrix dimension  $\max\{n_1, n_2\}$ , its dependency on other salient parameters — e.g. the number  $K$  of components, the ranks  $\{r_k\}$  of the ground-truth matrices  $\{\mathbf{M}_k^*\}$  — is likely sub-optimal. Improving the sample efficiency in these aspects is certainly an interesting direction to explore. In addition, in the presence of *random* noise, the performance of ScaledTGD saturates after the number of samples exceeds a certain threshold. It would be helpful to investigate other algorithms like expectation-maximization to see whether there is any performance gain one can harvest. Furthermore, our current theory builds upon the Gaussian designs  $\{\mathbf{A}_i\}$ , which often does not capture the practical scenarios. It is of great practical importance to develop efficient algorithms that can accommodate a wider range of design matrices  $\{\mathbf{A}_i\}$  — for instance, the case of mixed low-rank matrix completion. Last but not least, it would be of interest to study more general meta-learning settings in the presence of both light and heavy tasks (beyond the current single-sample setting) [KSS<sup>+</sup>20], and see how sample complexities can be reduced (compared to meta-learning for mixed regression) by exploiting such low-complexity structural priors.

## Acknowledgements

Y. Chen is supported in part by the grants AFOSR YIP award FA9550-19-1-0030, ONR N00014-19-1-2120, ARO YIP award W911NF-20-1-0097, ARO W911NF-18-1-0303, NSF CCF-1907661, IIS-1900140 and DMS-2014279, and the Princeton SEAS Innovation Award. H. V. Poor is supported in part by NSF CCF-1908308,

---

**Algorithm 5:** The tensor method for mixed linear regression [YCS16, Algorithm 1]

---

- 1 **Input:**  $\{\mathbf{a}_i, y_i\}_{1 \leq i \leq N}$ .
  - 2 Randomly split the samples into two disjoint sets  $\{\mathbf{a}_i, y_i\}_{1 \leq i \leq N_1}$ ,  $\{\mathbf{a}'_i, y'_i\}_{1 \leq i \leq N_2}$  such that  $N = N_1 + N_2$ , by assigning each sample to either dataset with probability 0.5.
  - 3 Compute  $m_0 \leftarrow \frac{1}{N_1} \sum_{i=1}^{N_1} y_i^2$ ,  $\mathbf{m}_1 \leftarrow \frac{1}{6N_2} \sum_{i=1}^{N_2} y_i'^3 \mathbf{a}'_i$ .
  - 4 Compute  $\mathbf{M}_2 \leftarrow \frac{1}{2N_1} \sum_{i=1}^{N_1} y_i^2 \mathbf{a}_i \mathbf{a}_i^\top - \frac{1}{2} m_0 \mathbf{I}_d$ ,  $\mathbf{M}_3 \leftarrow \frac{1}{6N_2} \sum_{i=1}^{N_2} y_i'^3 \mathbf{a}'_i \otimes^3 - \mathcal{T}(\mathbf{m}_1)$ , where  $\mathcal{T}$  is defined in (36).
  - 5 Denote the rank- $K$  SVD of  $\mathbf{M}_2$  as  $\mathbf{U}_2 \boldsymbol{\Sigma}_2 \mathbf{V}_2^\top$ , and compute the whitening matrix  $\mathbf{W} \leftarrow \mathbf{U}_2 \boldsymbol{\Sigma}_2^{-1/2}$ .
  - 6 Compute  $\widetilde{\mathbf{M}}_3 \leftarrow \mathbf{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})$ .
  - 7 Run the robust tensor power method [YCS16, Algorithm 2] on  $\widetilde{\mathbf{M}}_3$  to obtain  $K$  eigenvalue/eigenvector pairs  $\{\widetilde{\omega}_k, \widetilde{\boldsymbol{\beta}}_k\}_{1 \leq k \leq K}$ .
  - 8 Compute  $\omega_k \leftarrow 1/\widetilde{\omega}_k^2$ ,  $\boldsymbol{\beta}_k \leftarrow \widetilde{\omega}_k \mathbf{W}(\mathbf{W}^\top \mathbf{W})^{-1} \widetilde{\boldsymbol{\beta}}_k$ ,  $1 \leq k \leq K$ .
  - 9 **Output:**  $\{\omega_k, \boldsymbol{\beta}_k\}_{1 \leq k \leq K}$ .
- 

and in part by a Princeton Schmidt Data-X Research Award. We would like to thank Qixing Huang who taught us the symmetry synchronoziation problem in computer vision that largely inspired this research.

## A The tensor method for mixed linear regression

This section reviews the tensor method proposed in [YCS16] for solving mixed linear regression. For simplicity of exposition, we consider the noiseless case where we have access to the samples  $\{\mathbf{a}_i, y_i\}_{1 \leq i \leq N}$  obeying

$$y_i = \begin{cases} \langle \mathbf{a}_i, \boldsymbol{\beta}_1^* \rangle, & \text{if } i \in \Omega_1^*, \\ \dots & \dots \\ \langle \mathbf{a}_i, \boldsymbol{\beta}_K^* \rangle, & \text{if } i \in \Omega_K^*. \end{cases} \quad (35)$$

Our goal is to recover the ground truths  $\boldsymbol{\beta}_k^* \in \mathbb{R}^d$ ,  $1 \leq k \leq K$ , without knowing the index sets  $\{\Omega_k^*\}$ .

**Notation for tensors.** For two matrices  $\mathbf{A}$  and  $\mathbf{B}$ , denote by  $\mathbf{A} \otimes \mathbf{B}$  their Kronecker product, and let  $\mathbf{A} \otimes^3$  represent  $\mathbf{A} \otimes \mathbf{A} \otimes \mathbf{A}$ . For a symmetric tensor  $\mathbf{T} \in \mathbb{R}^{d \times d \times d}$  and matrices  $\mathbf{A} \in \mathbb{R}^{d \times d_1}$ ,  $\mathbf{B} \in \mathbb{R}^{d \times d_2}$ ,  $\mathbf{C} \in \mathbb{R}^{d \times d_3}$ , let  $\mathbf{T}(\mathbf{A}, \mathbf{B}, \mathbf{C}) \in \mathbb{R}^{d_1 \times d_2 \times d_3}$  denote the multi-linear matrix multiplication such that

$$[\mathbf{T}(\mathbf{A}, \mathbf{B}, \mathbf{C})]_{m,n,p} = \sum_{1 \leq i,j,k \leq d} T_{i,j,k} A_{i,m} B_{j,n} C_{k,p}, \quad 1 \leq m \leq d_1, 1 \leq n \leq d_2, 1 \leq p \leq d_3.$$

In addition, let  $\|\mathbf{T}\|$  stand for the operator norm of  $\mathbf{T}$ , namely,  $\|\mathbf{T}\| := \sup_{\mathbf{x}: \|\mathbf{x}\|_2=1} |\mathbf{T}(\mathbf{x}, \mathbf{x}, \mathbf{x})|$ .

**The tensor method: algorithm and rationale.** We summarize the tensor method in Algorithm 5, which is mostly the same as [YCS16, Algorithm 1] and included here for completeness.

In the following, we explain the intuitions behind its algorithmic design. Given data  $\{\mathbf{a}_i, y_i\}_{1 \leq i \leq N}$  generated according to (35), we compute the following empirical moments:

$$m_0 := \frac{1}{N} \sum_{i=1}^N y_i^2 \in \mathbb{R}, \quad \mathbf{m}_1 := \frac{1}{6N} \sum_{i=1}^N y_i^3 \mathbf{a}_i \in \mathbb{R}^d,$$

$$\mathbf{M}_2 := \frac{1}{2N} \sum_{i=1}^N y_i^2 \mathbf{a}_i \mathbf{a}_i^\top - \frac{1}{2} m_0 \mathbf{I}_d \in \mathbb{R}^{d \times d}, \quad \mathbf{M}_3 := \frac{1}{6N} \sum_{i=1}^N y_i^3 \mathbf{a}_i \otimes^3 - \mathcal{T}(\mathbf{m}_1) \in \mathbb{R}^{d \times d \times d},$$

here, letting  $\{\mathbf{e}_i\}_{1 \leq i \leq d}$  be the canonical basis of  $\mathbb{R}^d$ , we define the operator  $\mathcal{T}(\cdot) : \mathbb{R}^d \mapsto \mathbb{R}^{d \times d \times d}$  as

$$\mathcal{T}(\mathbf{m}) := \sum_{i=1}^d (\mathbf{m} \otimes \mathbf{e}_i \otimes \mathbf{e}_i + \mathbf{e}_i \otimes \mathbf{m} \otimes \mathbf{e}_i + \mathbf{e}_i \otimes \mathbf{e}_i \otimes \mathbf{m}), \quad \text{where } \mathbf{m} \in \mathbb{R}^d. \quad (36)$$

The key observation is that: under the Gaussian design (i.e.  $\mathbf{a}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ ),  $\mathbf{M}_2$  and  $\mathbf{M}_3$  reveal crucial second-order and third-order moments of  $\{\boldsymbol{\beta}_k^*\}$  since (cf. [YCS16, Lemma 1])

$$\mathbb{E}[\mathbf{M}_2] = \sum_{k=1}^K p_k \boldsymbol{\beta}_k^* (\boldsymbol{\beta}_k^*)^\top \quad \text{and} \quad \mathbb{E}[\mathbf{M}_3] = \sum_{k=1}^K p_k (\boldsymbol{\beta}_k^*)^{\otimes 3},$$

where we recall  $p_k = |\Omega_k^*|/N$ . This motivates one to apply tensor decomposition [AGH<sup>+</sup>14] on  $\mathbf{M}_2$  and  $\mathbf{M}_3$  in order to estimate  $\{\boldsymbol{\beta}_k^*\}$  and  $\{p_k\}$ . Indeed, the estimates  $\{\hat{\boldsymbol{\beta}}_k\}$  and  $\{\hat{\omega}_k\}$  returned by Algorithm 5 serve as our estimates of  $\{\boldsymbol{\beta}_k^*\}$  and  $\{p_k\}$ , respectively.

*Remark 1* (Sample splitting). Similar to [YCS16], we assume that  $m_0$  and  $\mathbf{M}_2$  are computed using one set of data, while  $\mathbf{M}_1$  and  $\mathbf{M}_3$  are obtained based on another *independent* set of samples. This sample splitting strategy ensures that the whitening matrix  $\mathbf{W}$  is independent of  $\mathbf{M}_3$ , thus simplifying theoretical analysis.

## B Proofs for Section 5

For notational simplicity, we use  $\text{dist}_{\mathcal{U}, \mathcal{V}}$  throughout to denote the following subspace estimation error:

$$\text{dist}_{\mathcal{U}, \mathcal{V}} := \max \left\{ \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}\|, \|\mathbf{V}\mathbf{V}^\top - \mathbf{V}^* \mathbf{V}^{*\top}\| \right\}. \quad (37)$$

### B.1 Proof of Theorem 3

The proof is decomposed into two steps: we first develop an upper bound  $\|\mathbf{Y} - \mathbb{E}[\mathbf{Y}]\|$  (where  $\mathbf{Y}$  is as defined in Algorithm 1), and then combine this with Wedin's Theorem to control the subspace distance  $\text{dist}_{\mathcal{U}, \mathcal{V}}$ .

**Step 1: controlling  $\|\mathbf{Y} - \mathbb{E}[\mathbf{Y}]\|$ .** We start by decomposing  $\mathbf{Y}$  into  $\mathbf{Y} = \mathbf{Y}_A + \mathbf{Y}_\zeta$ , where we define

$$\mathbf{Y}_A := \sum_{k=1}^K \frac{p_k}{|\Omega_k^*|} \sum_{i \in \Omega_k^*} \langle \mathbf{A}_i, \mathbf{M}_k^* \rangle \mathbf{A}_i \quad \text{and} \quad \mathbf{Y}_\zeta := \frac{1}{N} \sum_{i=1}^N \zeta_i \mathbf{A}_i.$$

Lemma 1 asserts that: with probability at least  $1 - Ce^{-cn}$  for some universal constants  $C, c > 0$ , we have

$$\left\| \frac{1}{|\Omega_k^*|} \sum_{i \in \Omega_k^*} \langle \mathbf{A}_i, \mathbf{M}_k^* \rangle \mathbf{A}_i - \mathbf{M}_k^* \right\| \leq \delta \|\mathbf{M}_k^*\|_{\text{F}}, \quad 1 \leq k \leq K,$$

as long as the sample size  $N$  satisfies (29), which together with the triangle inequality further implies

$$\|\mathbf{Y}_A - \mathbb{E}[\mathbf{Y}_A]\| \leq \sum_{k=1}^K p_k \left\| \frac{1}{|\Omega_k^*|} \sum_{i \in \Omega_k^*} \langle \mathbf{A}_i, \mathbf{M}_k^* \rangle \mathbf{A}_i - \mathbf{M}_k^* \right\| \leq \delta \sum_{k=1}^K p_k \|\mathbf{M}_k^*\|_{\text{F}} \leq \delta \max_{1 \leq k \leq K} \|\mathbf{M}_k^*\|_{\text{F}}.$$

In addition, [CP11, Lemma 1.1] reveals that with probability at least  $1 - Ce^{-cn}$  for some constants  $C, c > 0$ ,

$$\|\mathbf{Y}_\zeta\| \lesssim \sigma \sqrt{\frac{n}{N}} \lesssim \sigma \frac{\delta}{\sqrt{Kr}}$$

holds under the sample size condition (29). Given that  $\mathbb{E}[\mathbf{Y}_A] = \mathbb{E}[\mathbf{Y}] = \sum_k p_k \mathbf{M}_k^*$ , we have established the existence of some universal constant  $C_1 > 0$  such that

$$\|\mathbf{Y} - \mathbb{E}[\mathbf{Y}]\| \leq \|\mathbf{Y}_A - \mathbb{E}[\mathbf{Y}_A]\| + \|\mathbf{Y}_\zeta\| \leq C_1 \delta \left( \max_{1 \leq k \leq K} \|\mathbf{M}_k^*\|_{\text{F}} + \frac{\sigma}{\sqrt{Kr}} \right) =: \Delta. \quad (38)$$

**Step 2: controlling  $\text{dist}_{\mathbf{U}, \mathbf{V}}$ .** Before embarking on controlling  $\text{dist}_{\mathbf{U}, \mathbf{V}}$ , we make the following claim.

**Claim 1.** *Under the assumptions of Theorem 3, we have*

$$\text{col} \left\{ \sum_{k=1}^K p_k \mathbf{M}_k^* \right\} = \text{col} \left\{ [\mathbf{U}_1^*, \dots, \mathbf{U}_K^*] \right\}, \quad \text{row} \left\{ \sum_{k=1}^K p_k \mathbf{M}_k^* \right\} = \text{col} \left\{ [\mathbf{V}_1^*, \dots, \mathbf{V}_K^*] \right\}, \quad R = \sum_{k=1}^K r_k, \quad (39a)$$

$$\text{and} \quad \sigma_R \left( \sum_{k=1}^K p_k \mathbf{M}_k^* \right) \gtrsim \frac{1}{K} \min_k \sigma_{r_k}(\mathbf{M}_k^*). \quad (39b)$$

With this claim in place, we are ready to apply Wedin's Theorem [Wed72] to obtain

$$\text{dist}_{\mathbf{U}, \mathbf{V}} \leq \frac{\|\mathbf{Y} - \mathbb{E}[\mathbf{Y}]\|}{\sigma_R(\mathbb{E}[\mathbf{Y}]) - \|\mathbf{Y} - \mathbb{E}[\mathbf{Y}]\|} \leq \frac{\Delta}{\sigma_R(\mathbb{E}[\mathbf{Y}]) - \Delta} \leq \frac{2\Delta}{\sigma_R(\mathbb{E}[\mathbf{Y}])} = 2C_1 \frac{\delta \left( \max_k \|\mathbf{M}_k^*\|_{\text{F}} + \frac{\sigma}{\sqrt{Kr}} \right)}{\sigma_R \left( \sum_k p_k \mathbf{M}_k^* \right)}, \quad (40)$$

with the proviso that  $\Delta$  defined in (38) obeys  $\Delta \leq \frac{1}{2} \sigma_R(\mathbb{E}[\mathbf{Y}])$ . On the other hand, if instead one has  $\Delta > \frac{1}{2} \sigma_R(\mathbb{E}[\mathbf{Y}])$ , then we claim that (40) trivially holds; this can be seen by observing that  $\text{dist}_{\mathbf{U}, \mathbf{V}} \leq 1$ , while the right-hand side of (40) is greater than 1 if  $\Delta > \frac{1}{2} \sigma_R(\mathbb{E}[\mathbf{Y}])$ . Finally, Claim 1 tells us that

$$\sigma_R \left( \sum_k p_k \mathbf{M}_k^* \right) \gtrsim \frac{1}{K} \min_k \sigma_{r_k}(\mathbf{M}_k^*) \gtrsim \frac{1}{K \sqrt{r_K}} \min_k \|\mathbf{M}_k^*\|_{\text{F}}.$$

Substituting this relation into (40) immediately leads to the advertised bound (28) in Theorem 3.

*Proof of Claim 1.* Recall that we can write  $\sum_k p_k \mathbf{M}_k^*$  in terms of  $\{\mathbf{U}_k^*, \boldsymbol{\Sigma}_k^*, \mathbf{V}_k^*\}$ , in the form of (8). Therefore, to prove (39a), it suffices to show that  $\min\{\sigma_{R'}([\mathbf{U}_1^*, \dots, \mathbf{U}_K^*]), \sigma_{R'}([\mathbf{V}_1^*, \dots, \mathbf{V}_K^*])\} \geq 1/\sqrt{2}$ , where  $R' := \sum_k r_k$ . We only prove this for  $\sigma_{R'}([\mathbf{U}_1^*, \dots, \mathbf{U}_K^*])$ , since the proof for  $\sigma_{R'}([\mathbf{V}_1^*, \dots, \mathbf{V}_K^*])$  is identical. Denoting  $\mathbf{W} := [\mathbf{U}_1^*, \dots, \mathbf{U}_K^*]$  for notational convenience, we have

$$\mathbf{W}^\top \mathbf{W} = \begin{bmatrix} \mathbf{U}_1^{*\top} \\ \vdots \\ \mathbf{U}_K^{*\top} \end{bmatrix} [\mathbf{U}_1^* \quad \dots \quad \mathbf{U}_K^*] = \begin{bmatrix} \mathbf{I}_{r_1} & \mathbf{U}_1^{*\top} \mathbf{U}_2^* & \dots & \mathbf{U}_1^{*\top} \mathbf{U}_K^* \\ \mathbf{U}_2^{*\top} \mathbf{U}_1^* & \mathbf{I}_{r_2} & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \mathbf{U}_K^{*\top} \mathbf{U}_1^* & \dots & \dots & \mathbf{I}_{r_K} \end{bmatrix}.$$

This together with Assumption 1 gives

$$\|\mathbf{W}^\top \mathbf{W} - \mathbf{I}_{R'}\|_{\text{F}}^2 = \sum_{i \neq j} \|\mathbf{U}_i^{*\top} \mathbf{U}_j^*\|_{\text{F}}^2 \leq K^2 \left( \frac{1}{2K} \right)^2 \leq \frac{1}{4}.$$

Apply Weyl's inequality to obtain

$$\sigma_{R'}(\mathbf{W}^\top \mathbf{W}) \geq 1 - \|\mathbf{W}^\top \mathbf{W} - \mathbf{I}_{R'}\| \geq 1 - \|\mathbf{W}^\top \mathbf{W} - \mathbf{I}_{R'}\|_{\text{F}} \geq \frac{1}{2},$$

thus indicating that  $\sigma_{R'}(\mathbf{W}) = \sqrt{\sigma_{R'}(\mathbf{W}^\top \mathbf{W})} \geq 1/\sqrt{2}$ . This completes the proof of (39a).

Next, we turn attention to (39b). Denote the SVD of  $[\mathbf{U}_1^*, \dots, \mathbf{U}_K^*]$  (resp.  $[\mathbf{V}_1^*, \dots, \mathbf{V}_K^*]$ ) as  $\mathbf{U}_{\text{left}} \boldsymbol{\Sigma}_{\text{left}} \mathbf{V}_{\text{left}}^\top$  (resp.  $\mathbf{U}_{\text{right}} \boldsymbol{\Sigma}_{\text{right}} \mathbf{V}_{\text{right}}^\top$ ), where  $\mathbf{V}_{\text{left}}$  (resp.  $\mathbf{V}_{\text{right}}$ ) is a  $R \times R$  orthonormal matrix. Substitution into (8) yields

$$\sum_{k=1}^K p_k \mathbf{M}_k^* = \mathbf{U}_{\text{left}} \boldsymbol{\Sigma}_{\text{left}} \mathbf{V}_{\text{left}}^\top \text{diag}(\{p_k \boldsymbol{\Sigma}_k^*\}_{1 \leq k \leq K}) \mathbf{V}_{\text{right}} \boldsymbol{\Sigma}_{\text{right}} \mathbf{U}_{\text{right}}^\top,$$

where  $\text{diag}(\{p_k \boldsymbol{\Sigma}_k^*\}_{1 \leq k \leq K})$  is a  $R \times R$  full-rank diagonal matrix, with blocks  $p_1 \boldsymbol{\Sigma}_1^*, \dots, p_K \boldsymbol{\Sigma}_K^*$  on the diagonal. This implies that

$$\sigma_R \left( \sum_{k=1}^K p_k \mathbf{M}_k^* \right) = \sigma_R \left( \boldsymbol{\Sigma}_{\text{left}} \mathbf{V}_{\text{left}}^\top \text{diag}(\{p_k \boldsymbol{\Sigma}_k^*\}_{1 \leq k \leq K}) \mathbf{V}_{\text{right}} \boldsymbol{\Sigma}_{\text{right}} \right) \geq \sigma_R(\boldsymbol{\Sigma}_{\text{left}}) \sigma_R(\boldsymbol{\Sigma}_{\text{right}}) \cdot \min_k \{p_k \sigma_{r_k}(\mathbf{M}_k^*)\}$$

$$\geq \left(\frac{1}{\sqrt{2}}\right)^2 \min_k \{p_k \sigma_{r_k}(\mathbf{M}_k^*)\} \gtrsim \frac{1}{K} \min_k \sigma_{r_k}(\mathbf{M}_k^*),$$

where the last inequality uses the assumption that  $p_k \gtrsim 1/K$ . This establishes (39b).  $\square$

## B.2 Proof of Theorem 4

**Step 1: basic properties of the auxiliary mixed linear regression problem.** We begin by formally characterizing the intermediate mixed linear regression problem in Stage 2. It is easily seen from Section 2.2 that for  $i \in \Omega_k^*$ , one has

$$y_i = \langle \mathbf{A}_i, \mathbf{M}_k^* \rangle + \zeta_i = \langle \mathbf{a}_i, \boldsymbol{\beta}_k \rangle + \underbrace{z_i + \zeta_i}_{=:\xi_i}, \quad (41)$$

where the additional term

$$z_i := \langle \mathbf{A}_i, \mathbf{M}_k^* \rangle - \langle \mathbf{a}_i, \boldsymbol{\beta}_k \rangle = \langle \mathbf{A}_i, \mathbf{M}_k^* - \mathbf{U}\mathbf{U}^\top \mathbf{M}_k^* \mathbf{V}\mathbf{V}^\top \rangle \quad (42)$$

accounts for the subspace estimation error. In words, the observations  $\{y_i\}$  can be equivalently written in the mixed linear regression form, where  $\{\boldsymbol{\beta}_k\}$  constitutes the underlying parameters,  $\{\mathbf{a}_i\}$  the measurement vectors and  $\{\xi_i\}$  the measurement noise. We then focus on characterizing the properties of  $\mathbf{a}_i$  and  $\xi_i$ .

Recall from Algorithm 2 that  $\mathbf{a}_i = \text{vec}(\mathbf{U}^\top \mathbf{A}_i \mathbf{V})$ . In view of the independence between  $\{\mathbf{A}_i\}$  and  $\mathbf{U}, \mathbf{V}$ , one can deduce that

$$\mathbf{a}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d), \quad 1 \leq i \leq N,$$

where  $d := R^2$ . Again, leveraging the independence between  $\{\mathbf{A}_i, \zeta_i\}$  and  $\mathbf{U}, \mathbf{V}$ , we have

$$\xi_i = \langle \mathbf{A}_i, \mathbf{M}_k^* - \mathbf{U}\mathbf{U}^\top \mathbf{M}_k^* \mathbf{V}\mathbf{V}^\top \rangle + \zeta_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \|\mathbf{M}_k^* - \mathbf{U}\mathbf{U}^\top \mathbf{M}_k^* \mathbf{V}\mathbf{V}^\top\|_{\mathbb{F}}^2 + \sigma^2).$$

For notational convenience, we shall denote the variance to be

$$\sigma_k^2 := \|\mathbf{M}_k^* - \mathbf{U}\mathbf{U}^\top \mathbf{M}_k^* \mathbf{V}\mathbf{V}^\top\|_{\mathbb{F}}^2 + \sigma^2, \quad 1 \leq k \leq K. \quad (43)$$

More importantly, the measurement vectors  $\{\mathbf{a}_i\}$  are independent of the measurement noise  $\{\xi_i\}$ . To see this, one has

$$\begin{aligned} \mathbb{E}[\xi_i \mathbf{a}_i] &= \mathbb{E}[\zeta_i \mathbf{a}_i] + \mathbb{E}[z_i \mathbf{a}_i] = \mathbf{0} + \text{vec}(\mathbb{E}[\langle \mathbf{A}_i, \mathbf{M}_k^* - \mathbf{U}\mathbf{U}^\top \mathbf{M}_k^* \mathbf{V}\mathbf{V}^\top \rangle \mathbf{U}^\top \mathbf{A}_i \mathbf{V}]) \\ &= \text{vec}(\mathbf{U}^\top (\mathbf{M}_k^* - \mathbf{U}\mathbf{U}^\top \mathbf{M}_k^* \mathbf{V}\mathbf{V}^\top) \mathbf{V}) = \mathbf{0}. \end{aligned}$$

Here the second equality follows from the independence between  $\zeta_i$  and  $\mathbf{A}_i, \mathbf{U}, \mathbf{V}$ , whereas the last line utilizes the independence between  $\mathbf{A}_i$  and  $\mathbf{U}, \mathbf{V}$  and the isotropic property of  $\mathbf{A}_i$ .

In conclusion, in Line 3 of Algorithm 2, we are equivalently faced with a  $d$ -dimensional mixed linear regression problem with data  $\{\mathbf{a}_i, y_i\}_{1 \leq i \leq N}$ , which satisfies that for  $i \in \Omega_k^*$ ,

$$y_i = \langle \mathbf{a}_i, \boldsymbol{\beta}_k \rangle + \xi_i, \quad \xi_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_k^2), \quad \mathbf{a}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d) \quad (44)$$

with  $\xi_i$  being independent from  $\mathbf{a}_i$ .

**Step 2: performance of the tensor method.** Next, we characterize the performance of the tensor method for solving the above mixed linear regression problem. Our proof follows closely that of [YCS16, Theorem 1], with minor modifications to accommodate the noise  $\{\xi_i\}$ . Therefore we only provide a sketch here.

Recall that in Algorithm 5, we randomly split the input data  $\{\mathbf{a}_i, y_i\}_{1 \leq i \leq N}$  into two sets  $\{\mathbf{a}_i, y_i\}_{1 \leq i \leq N_1}$  and  $\{\mathbf{a}'_i, y'_i\}_{1 \leq i \leq N_2}$  (with slight abuse of notation). This sample splitting strategy is adopted merely to decouple statistical dependence and facilitate analysis. The high-level idea of the proof of [YCS16, Theorem 1] is simple to state: if the quantities

$$\left\| \mathbf{M}_2 - \sum_{k=1}^K p_k \boldsymbol{\beta}_k \boldsymbol{\beta}_k^\top \right\| \quad \text{and} \quad \left\| \left( \mathbf{M}_3 - \sum_{k=1}^K p_k \boldsymbol{\beta}_k^{\otimes 3} \right) (\mathbf{W}, \mathbf{W}, \mathbf{W}) \right\| \quad (45)$$



are sufficiently small, then the tensor method returns reliable estimates of  $\{\beta_k\}$ ; see [YCS16, Eq. (24) in Section 5.4.1]. Here, the empirical moments  $\mathbf{M}_2, \mathbf{M}_3$  and the whitening matrix  $\mathbf{W}$  are defined in Algorithm 5.

With this connection in place, it suffices to control the quantities in (45). While the analysis in [YCS16, Section 5.4.2] only applies to the noiseless mixed linear regression problem, we can easily modify it to accommodate our noisy case (44). The trick is to *augment*  $\{\beta_k\}$  and  $\{\mathbf{a}_i\}$  as follows:

$$\beta_k^{\text{aug}} := \begin{bmatrix} \beta_k \\ \sigma_k \end{bmatrix} \in \mathbb{R}^{d+1}, \quad 1 \leq k \leq K; \quad \mathbf{a}_i^{\text{aug}} := \begin{bmatrix} \mathbf{a}_i \\ \xi_i/\sigma_k \end{bmatrix} \in \mathbb{R}^{d+1}, \quad i \in \Omega_k^*. \quad (46)$$

The advantage is clear: the noisy mixed linear regression problem (44) can be equivalently phrased as a noiseless one, that is for all  $i \in \Omega_k^*$ ,

$$\mathbf{a}_i^{\text{aug}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_{d+1}) \quad \text{and} \quad y_i = \langle \mathbf{a}_i^{\text{aug}}, \beta_k^{\text{aug}} \rangle. \quad (47)$$

Similarly, we can define  $\mathbf{a}_i^{\text{aug}'}$  analogously, and introduce the augmented versions of the empirical moments as follows:

$$m_0^{\text{aug}} := \frac{1}{N_1} \sum_{i=1}^{N_1} y_i^2 \in \mathbb{R}, \quad \mathbf{m}_1^{\text{aug}} := \frac{1}{6N_2} \sum_{i=1}^{N_2} y_i'^3 \mathbf{a}_i^{\text{aug}'} \in \mathbb{R}^{d+1}, \quad (48a)$$

$$\mathbf{M}_2^{\text{aug}} := \frac{1}{2N_1} \sum_{i=1}^{N_1} y_i^2 \mathbf{a}_i^{\text{aug}} (\mathbf{a}_i^{\text{aug}})^\top - \frac{1}{2} m_0^{\text{aug}} \mathbf{I}_{d+1} \in \mathbb{R}^{(d+1) \times (d+1)}, \quad (48b)$$

$$\mathbf{M}_3^{\text{aug}} := \frac{1}{6N_2} \sum_{i=1}^{N_2} y_i'^3 (\mathbf{a}_i^{\text{aug}'})^{\otimes 3} - \mathcal{T}^{\text{aug}}(\mathbf{m}_1^{\text{aug}}) \in \mathbb{R}^{(d+1) \times (d+1) \times (d+1)}, \quad (48c)$$

where  $\mathcal{T}^{\text{aug}}(\cdot)$  is defined analogously as in (36). By virtue of the augmentation procedure,  $\mathbf{M}_2$  (resp.  $\mathbf{M}_3$ ) is a sub-matrix (resp. sub-tensor) of  $\mathbf{M}_2^{\text{aug}}$  (resp.  $\mathbf{M}_3^{\text{aug}}$ ). Consequently, we have

$$\begin{aligned} \left\| \mathbf{M}_2 - \sum_{k=1}^K p_k \beta_k \beta_k^\top \right\| &\leq \left\| \mathbf{M}_2^{\text{aug}} - \sum_{k=1}^K p_k \beta_k^{\text{aug}} (\beta_k^{\text{aug}})^\top \right\|; \\ \left\| \left( \mathbf{M}_3 - \sum_{k=1}^K p_k \beta_k^{\otimes 3} \right) (\mathbf{W}, \mathbf{W}, \mathbf{W}) \right\| &= \left\| \left( \mathbf{M}_3^{\text{aug}} - \sum_{k=1}^K p_k (\beta_k^{\text{aug}})^{\otimes 3} \right) (\mathbf{W}^{\text{aug}}, \mathbf{W}^{\text{aug}}, \mathbf{W}^{\text{aug}}) \right\|, \end{aligned}$$

where  $\mathbf{W}^{\text{aug}} := [\mathbf{W}^\top, \mathbf{0}]^\top$ .

With the above augmented vectors/matrices/tensors in place, one can follow the analysis in [YCS16, Section 5.4.2] to upper bound the quantities above. One subtle issue is that our sampling scheme is slightly different from the one in [YCS16], where each sample has i.i.d. labeling; nevertheless, it is easy to check that this difference is minor, and does not affect the result of the analysis. Indeed, repeating the analysis in [YCS16, Section 5.4] yields the conclusion that: in order to achieve  $\epsilon$  errors (30) with probability at least  $1 - \gamma$ , it suffices to require the sample complexities to exceed (analogous to [YCS16, Eq. (13)])

$$N_1 \geq C_1 \left( \frac{d}{(\min_k p_k) \epsilon^2} \frac{\max_k \|\beta_k^{\text{aug}}\|_2^{10}}{\sigma_K (\sum_k p_k \beta_k \beta_k^\top)^5} \log \frac{12K}{\gamma} \log^2 N_1 + \frac{K}{(\min_k p_k) \gamma} \right) \quad (49a)$$

$$\stackrel{\text{(i)}}{\asymp} \frac{dK^6}{\epsilon^2} \left( \Gamma^{10} + \frac{\sigma^{10}}{\min_k \|\mathbf{M}_k^*\|_{\text{F}}^{10}} \right) \log(K \log n) \log^2 N_1 + K^2 \log n, \quad (49b)$$

$$N_2 \geq C_2 \left( \frac{(K^2 + d)}{(\min_k p_k) \epsilon^2} \frac{\max_k \|\beta_k^{\text{aug}}\|_2^6}{\sigma_K (\sum_k p_k \beta_k \beta_k^\top)^3} \log \frac{12K}{\gamma} \log^3 N_2 + \frac{K}{(\min_k p_k) \gamma} \right) \quad (49c)$$

$$\stackrel{\text{(ii)}}{\asymp} \frac{dK^4}{\epsilon^2} \left( \Gamma^6 + \frac{\sigma^6}{\min_k \|\mathbf{M}_k^*\|_{\text{F}}^6} \right) \log(K \log n) \log^3 N_2 + K^2 \log n. \quad (49d)$$

Here  $C_1, C_2 > 0$  are some sufficiently large constants, and the simplifications (i) (ii) hold due to the following facts: (i)  $d = R^2 \geq K^2$ , (ii)  $\min_k p_k \asymp 1/K$ , (iii) we choose  $\gamma = O(1/\log n)$ , (iv)  $\|\beta_k^{\text{aug}}\|_2^2 = \|\beta_k\|_2^2 + \sigma_k^2 = \|\beta_k\|_2^2 + \|\mathbf{M}_k^* - \mathbf{U}\mathbf{U}^\top \mathbf{M}_k^* \mathbf{V}\mathbf{V}^\top\|_{\text{F}}^2 + \sigma^2$ , and (v) the following claim (in particular, (51) and (52) therein).

**Claim 2.** *Instate the assumptions of Theorem 4.*

1. The ground-truth matrices  $\{\mathbf{M}_k^*\}_{1 \leq k \leq K}$  satisfy that for all  $1 \leq i, j \leq K, i \neq j$ ,

$$|\langle \mathbf{M}_i^*, \mathbf{M}_j^* \rangle| \leq \frac{1}{4K\Gamma^2} \|\mathbf{M}_i^*\|_{\mathbb{F}} \|\mathbf{M}_j^*\|_{\mathbb{F}}, \quad \text{and} \quad \|\mathbf{M}_i^* - \mathbf{M}_j^*\|_{\mathbb{F}} \gtrsim \|\mathbf{M}_i^*\|_{\mathbb{F}} + \|\mathbf{M}_j^*\|_{\mathbb{F}}. \quad (50)$$

2. In addition, the parameters  $\{\boldsymbol{\beta}_k\}_{1 \leq k \leq K}$  obey that for all  $1 \leq k, i, j \leq K, i \neq j$ ,

$$0.9 \|\mathbf{M}_k^*\|_{\mathbb{F}} \leq \|\boldsymbol{\beta}_k\|_2 \leq \|\mathbf{M}_k^*\|_{\mathbb{F}}, \quad \text{and} \quad |\langle \boldsymbol{\beta}_i, \boldsymbol{\beta}_j \rangle| \leq \frac{1}{2K\Gamma^2} \|\boldsymbol{\beta}_i\|_2 \|\boldsymbol{\beta}_j\|_2. \quad (51)$$

3. In the end, we have

$$\sigma_K \left( \sum_{k=1}^K p_k \boldsymbol{\beta}_k \boldsymbol{\beta}_k^\top \right) \asymp \frac{1}{K} \min_{1 \leq k \leq K} \|\mathbf{M}_k^*\|_{\mathbb{F}}^2. \quad (52)$$

Armed with (49b) and (49d), we can plug in the bounds  $d = R^2 \leq K^2 r^2$  and  $\log(K \log n) \lesssim \log n$  to complete the proof of Theorem 4.

*Proof of Claim 2.* With regards to the first part of (50), it is seen that

$$\begin{aligned} |\langle \mathbf{M}_i^*, \mathbf{M}_j^* \rangle| &= |\langle \mathbf{U}_i^* \boldsymbol{\Sigma}_i^* \mathbf{V}_i^{*\top}, \mathbf{U}_j^* \boldsymbol{\Sigma}_j^* \mathbf{V}_j^{*\top} \rangle| = |\langle \boldsymbol{\Sigma}_i^*, \mathbf{U}_i^{*\top} \mathbf{U}_j^* \boldsymbol{\Sigma}_j^* \mathbf{V}_j^{*\top} \mathbf{V}_i^* \rangle| \leq \|\mathbf{U}_i^{*\top} \mathbf{U}_j^*\|_{\mathbb{F}} \|\mathbf{V}_i^{*\top} \mathbf{V}_j^*\|_{\mathbb{F}} \|\boldsymbol{\Sigma}_i^*\|_{\mathbb{F}} \|\boldsymbol{\Sigma}_j^*\|_{\mathbb{F}} \\ &\leq \left( \frac{1}{2\sqrt{K\Gamma}} \right)^2 \|\boldsymbol{\Sigma}_i^*\|_{\mathbb{F}} \|\boldsymbol{\Sigma}_j^*\|_{\mathbb{F}} = \frac{1}{4K\Gamma^2} \|\mathbf{M}_i^*\|_{\mathbb{F}} \|\mathbf{M}_j^*\|_{\mathbb{F}}, \end{aligned}$$

where the second line utilizes Assumption 1. The second part of (50) follows immediately from the first part and some elementary calculations.

Next, we turn to proving (51). Recall the definitions  $\boldsymbol{\beta}_k = \text{vec}(\mathbf{S}_k) = \text{vec}(\mathbf{U}^\top \mathbf{M}_k^* \mathbf{V})$  and  $\text{dist}_{\mathbf{U}, \mathbf{V}} = \max\{\|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}\|, \|\mathbf{V}\mathbf{V}^\top - \mathbf{V}^* \mathbf{V}^{*\top}\|\}$ . We have the upper bound  $\|\boldsymbol{\beta}_k\|_2 = \|\mathbf{U}^\top \mathbf{M}_k^* \mathbf{V}\|_{\mathbb{F}} \leq \|\mathbf{M}_k^*\|_{\mathbb{F}}$  as well as the lower bound

$$\begin{aligned} \|\boldsymbol{\beta}_k\|_2 &= \|\mathbf{U}^\top \mathbf{M}_k^* \mathbf{V}\|_{\mathbb{F}} = \|\mathbf{U}\mathbf{U}^\top \mathbf{M}_k^* \mathbf{V}\mathbf{V}^\top\|_{\mathbb{F}} \geq \|\mathbf{M}_k^*\|_{\mathbb{F}} - \|\mathbf{U}^* \mathbf{U}^{*\top} \mathbf{M}_k^* \mathbf{V}^* \mathbf{V}^{*\top} - \mathbf{U}\mathbf{U}^\top \mathbf{M}_k^* \mathbf{V}\mathbf{V}^\top\|_{\mathbb{F}} \\ &\geq \|\mathbf{M}_k^*\|_{\mathbb{F}} - \|(\mathbf{U}^* \mathbf{U}^{*\top} - \mathbf{U}\mathbf{U}^\top) \mathbf{M}_k^* \mathbf{V}^* \mathbf{V}^{*\top}\|_{\mathbb{F}} - \|\mathbf{U}\mathbf{U}^\top \mathbf{M}_k^* (\mathbf{V}^* \mathbf{V}^{*\top} - \mathbf{V}\mathbf{V}^\top)\|_{\mathbb{F}} \\ &\geq (1 - 2 \text{dist}_{\mathbf{U}, \mathbf{V}}) \|\mathbf{M}_k^*\|_{\mathbb{F}} \geq 0.9 \|\mathbf{M}_k^*\|_{\mathbb{F}}, \end{aligned}$$

where the last inequality uses the assumption that  $\text{dist}_{\mathbf{U}, \mathbf{V}} \leq c_1/(K\Gamma^2) \leq 0.05$ ; this justifies the first part of (51). To prove the second part of (51), we start with the decomposition

$$\begin{aligned} \langle \boldsymbol{\beta}_i, \boldsymbol{\beta}_j \rangle &= \langle \mathbf{U}^\top \mathbf{M}_i^* \mathbf{V}, \mathbf{U}^\top \mathbf{M}_j^* \mathbf{V} \rangle = \langle \mathbf{M}_i^*, \mathbf{U}\mathbf{U}^\top \mathbf{M}_j^* \mathbf{V}\mathbf{V}^\top \rangle \\ &= \langle \mathbf{M}_i^*, \mathbf{M}_j^* \rangle + \langle \mathbf{M}_i^*, \mathbf{U}\mathbf{U}^\top \mathbf{M}_j^* \mathbf{V}\mathbf{V}^\top - \mathbf{U}^* \mathbf{U}^{*\top} \mathbf{M}_j^* \mathbf{V}^* \mathbf{V}^{*\top} \rangle \\ &\quad + \langle \mathbf{M}_i^*, \mathbf{U}^* \mathbf{U}^{*\top} \mathbf{M}_j^* \mathbf{V}^* \mathbf{V}^{*\top} - \mathbf{U}^* \mathbf{U}^{*\top} \mathbf{M}_j^* \mathbf{V}^* \mathbf{V}^{*\top} \rangle, \end{aligned}$$

which together with the triangle inequality yields

$$|\langle \boldsymbol{\beta}_i, \boldsymbol{\beta}_j \rangle| \leq |\langle \mathbf{M}_i^*, \mathbf{M}_j^* \rangle| + 2 \text{dist}_{\mathbf{U}, \mathbf{V}} \|\mathbf{M}_i^*\|_{\mathbb{F}} \|\mathbf{M}_j^*\|_{\mathbb{F}}.$$

In light of the first part of (50), the first part of (51), and our assumption on  $\text{dist}_{\mathbf{U}, \mathbf{V}}$ , this establishes the second part of (51).

Finally, it remains to prove (52). In view of the assumption that  $p_k \asymp 1/K$  ( $1 \leq k \leq K$ ), one has

$$\sigma_K \left( \sum_{k=1}^K p_k \boldsymbol{\beta}_k \boldsymbol{\beta}_k^\top \right) \asymp \frac{1}{K} \sigma_K \left( \sum_{k=1}^K \boldsymbol{\beta}_k \boldsymbol{\beta}_k^\top \right). \quad (53)$$

Therefore, it suffices to show that  $\sigma_K(\sum_k \boldsymbol{\beta}_k \boldsymbol{\beta}_k^\top) \asymp \min_k \|\mathbf{M}_k^*\|_{\mathbb{F}}^2$ . Towards this, we find it helpful to define  $\mathbf{B} := [\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_K] \in \mathbb{R}^{d \times K}$ , and decompose  $\mathbf{B}^\top \mathbf{B}$  as  $\mathbf{B}^\top \mathbf{B} = \mathbf{D} + \mathbf{O}$ . Here,  $\mathbf{D}$  stands for the diagonal

part of  $\mathbf{B}^\top \mathbf{B}$  with  $D_{kk} = \|\beta_k\|_2^2$ , while  $\mathbf{O}$  is the off-diagonal part of  $\mathbf{B}^\top \mathbf{B}$ . Note that for any  $i \neq j$ ,  $[\mathbf{O}]_{ij} = [\mathbf{B}^\top \mathbf{B}]_{ij} = \langle \beta_i, \beta_j \rangle$ , which combined with (51) gives

$$\|\mathbf{O}\|_{\mathbb{F}}^2 = \sum_{i \neq j} \langle \beta_i, \beta_j \rangle^2 \leq K^2 \left( \frac{1}{2K\Gamma^2} \right)^2 \max_k \|\beta_k\|_2^4 = \frac{1}{4\Gamma^4} \max_k \|\beta_k\|_2^4 \leq \frac{1}{2} \min_k \|\beta_k\|_2^4;$$

the last inequality follows from the definition  $\Gamma = \max_k \|\mathbf{M}_k^*\|_{\mathbb{F}} / \min_k \|\mathbf{M}_k^*\|_{\mathbb{F}}$ , and the first part of (51). This together with Weyl's inequality implies that

$$\left| \sigma_K(\mathbf{B}^\top \mathbf{B}) - \sigma_K(\mathbf{D}) \right| = \left| \sigma_K(\mathbf{B}^\top \mathbf{B}) - \min_k \|\beta_k\|_2^2 \right| \leq \|\mathbf{O}\|_{\mathbb{F}} \leq \frac{1}{\sqrt{2}} \min_k \|\beta_k\|_2^2. \quad (54)$$

As a result, we arrive at

$$\sigma_K \left( \sum_k \beta_k \beta_k^\top \right) = \sigma_K(\mathbf{B}\mathbf{B}^\top) = \sigma_K(\mathbf{B}^\top \mathbf{B}) \asymp \min_k \|\beta_k\|_2^2 \asymp \min_k \|\mathbf{M}_k^*\|_{\mathbb{F}}^2,$$

which in conjunction with (53) completes the proof of (52).  $\square$

### B.3 Proof of Proposition 1

To begin with, the triangle inequality gives

$$\|\mathbf{L}_k \mathbf{R}_k^\top - \mathbf{M}_k^*\|_{\mathbb{F}} \leq \|\mathbf{U} \mathbf{S}_k \mathbf{V}^\top - \mathbf{M}_k^*\|_{\mathbb{F}} + \|\mathbf{L}_k \mathbf{R}_k^\top - \mathbf{U} \mathbf{S}_k \mathbf{V}^\top\|_{\mathbb{F}}. \quad (55)$$

Regarding the first term on the right-hand side of (55), we plug in the definition (13) of  $\mathbf{S}_k$  to obtain

$$\begin{aligned} \|\mathbf{U} \mathbf{S}_k \mathbf{V}^\top - \mathbf{M}_k^*\|_{\mathbb{F}} &= \|\mathbf{U} \mathbf{U}^\top \mathbf{M}_k^* \mathbf{V} \mathbf{V}^\top - \mathbf{U}^* \mathbf{U}^{*\top} \mathbf{M}_k^* \mathbf{V}^* \mathbf{V}^{*\top}\|_{\mathbb{F}} \\ &\leq \|(\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}) \mathbf{M}_k^* \mathbf{V} \mathbf{V}^\top\|_{\mathbb{F}} + \|\mathbf{U}^* \mathbf{U}^{*\top} \mathbf{M}_k^* (\mathbf{V} \mathbf{V}^\top - \mathbf{V}^* \mathbf{V}^{*\top})\|_{\mathbb{F}} \\ &\leq 2 \max \left\{ \|\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}\|, \|\mathbf{V} \mathbf{V}^\top - \mathbf{V}^* \mathbf{V}^{*\top}\| \right\} \|\mathbf{M}_k^*\|_{\mathbb{F}}. \end{aligned}$$

With regards to the second term on the right-hand side of (55), we observe that

$$\begin{aligned} \|\mathbf{L}_k \mathbf{R}_k^\top - \mathbf{U} \mathbf{S}_k \mathbf{V}^\top\|_{\mathbb{F}} &\leq \|\mathbf{L}_k \mathbf{R}_k^\top - \mathbf{U} \widehat{\mathbf{S}}_k \mathbf{V}^\top\|_{\mathbb{F}} + \|\mathbf{U} \widehat{\mathbf{S}}_k \mathbf{V}^\top - \mathbf{U} \mathbf{S}_k \mathbf{V}^\top\|_{\mathbb{F}} \\ &\stackrel{(i)}{\leq} 2 \|\mathbf{U} (\widehat{\mathbf{S}}_k - \mathbf{S}_k) \mathbf{V}^\top\|_{\mathbb{F}} \stackrel{(ii)}{\leq} 2 \|\widehat{\beta}_k - \beta_k\|_2. \end{aligned}$$

Here, (i) follows since  $\mathbf{L}_k \mathbf{R}_k^\top$  is the best rank- $r_k$  approximation of  $\mathbf{U} \widehat{\mathbf{S}}_k \mathbf{V}^\top$  and  $\mathbf{U} \mathbf{S}_k \mathbf{V}^\top$  is also rank- $r_k$ ; (ii) holds since  $\widehat{\beta}_k = \text{vec}(\widehat{\mathbf{S}}_k)$  and  $\beta_k = \text{vec}(\mathbf{S}_k)$ . Substitution into (55) establishes (32).

### B.4 Proof of Theorem 5

We shall only prove the local convergence w.r.t. the matrix  $\mathbf{M}_1^*$ ; the proof for other components is identical and hence is omitted. Our proof is decomposed into three steps.

1. Study the ScaledTGD dynamics (particularly the population-level dynamics), and control the effects of mislabeling and finite-sample errors.
2. Show that if the estimation error is larger than the error floor (namely, the last term in (34)), then one step of the ScaledTGD update contracts the error by a constant factor.
3. Show that, once the estimation error gets smaller than this error floor, then the estimation errors remain small in subsequent iterations.

Before continuing, we note that Condition (33) with  $k = 1$  implies the existence of some constant  $c_1 > 0$  such that

$$(\mathbf{L}^0, \mathbf{R}^0) \in \mathcal{B},$$

where

$$\mathcal{B} := \left\{ (\mathbf{L}, \mathbf{R}) \in \mathbb{R}^{n_1 \times r_1} \times \mathbb{R}^{n_2 \times r_1} : \|\mathbf{L}\mathbf{R}^\top - \mathbf{M}_1^*\|_{\text{F}} \leq c_1 \min \left\{ \sigma_{r_1}(\mathbf{M}_1^*), \frac{1}{K} \min_{j \neq 1} \|\mathbf{M}_j^* - \mathbf{M}_1^*\|_{\text{F}} \right\} \right\}. \quad (56)$$

This arises from the inequalities  $\sigma_{r_1}(\mathbf{M}_1^*) \geq \|\mathbf{M}_1^*\|_{\text{F}} / (\sqrt{r}\kappa)$  and  $\min_{j \neq 1} \|\mathbf{M}_j^* - \mathbf{M}_1^*\|_{\text{F}} \gtrsim \|\mathbf{M}_1^*\|_{\text{F}}$  (due to Assumption 1). We isolate Condition (56) since it is more convenient to work with in the analysis.

**Notation.** To simplify presentation, we shall often let  $(\mathbf{L}, \mathbf{R})$  denote an iterate lying within  $\mathcal{B}$  (cf. (56)), and define the corresponding estimation errors as

$$\Delta_k := \mathbf{L}\mathbf{R}^\top - \mathbf{M}_k^*, \quad 1 \leq k \leq K. \quad (57)$$

The truncating level for a prescribed truncating fraction  $\alpha$  is denoted by

$$\tau := Q_\alpha \left( \left\{ |\langle \mathbf{A}_i, \mathbf{L}\mathbf{R}^\top \rangle - y_i| \right\}_{1 \leq i \leq N} \right), \quad (58)$$

where  $Q_\alpha$  is the  $\alpha$ -quantile defined in Section 1.2. We also define the following functions and quantities:

$$\mathbb{1}(a; b) := \mathbb{1}(|a| \leq b), \quad a, b \in \mathbb{R}, \quad (59)$$

$$w(x) := \int_{-x}^x t^2 \phi(t) dt, \quad x \geq 0, \quad w_k := w \left( \frac{\tau}{\sqrt{\|\Delta_k\|_{\text{F}}^2 + \sigma^2}} \right), \quad 1 \leq k \leq K, \quad (60)$$

where  $\phi$  stands for the probability density function of a standard Gaussian random variable.

**Step 1: characterizing the ScaledTGD dynamic.** The above notation allows one to express the ScaledTGD update rule (16) as

$$\mathbf{L}^+ = \mathbf{L} - \frac{\eta}{N} \sum_{i=1}^N (\langle \mathbf{A}_i, \mathbf{L}\mathbf{R}^\top \rangle - y_i) \mathbb{1}(\langle \mathbf{A}_i, \mathbf{L}\mathbf{R}^\top \rangle - y_i; \tau) \mathbf{A}_i \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1}, \quad (61a)$$

$$\mathbf{R}^+ = \mathbf{R} - \frac{\eta}{N} \sum_{i=1}^N (\langle \mathbf{A}_i, \mathbf{L}\mathbf{R}^\top \rangle - y_i) \mathbb{1}(\langle \mathbf{A}_i, \mathbf{L}\mathbf{R}^\top \rangle - y_i; \tau) \mathbf{A}_i^\top \mathbf{L} (\mathbf{L}^\top \mathbf{L})^{-1}. \quad (61b)$$

Recall that for any  $i \in \Omega_k^*$ , we have  $y_i = \langle \mathbf{A}_i, \mathbf{M}_k^* \rangle + \zeta_i$ , and thus

$$\langle \mathbf{A}_i, \mathbf{L}\mathbf{R}^\top \rangle - y_i = \langle \mathbf{A}_i, \Delta_k \rangle - \zeta_i, \quad \text{for all } i \in \Omega_k^*. \quad (62)$$

The following result makes apparent a useful decomposition of the ScaledTGD update rule.

**Claim 3.** Recall the notation (59) and (60). The ScaledTGD update rule (61) can be written as

$$\mathbf{L}^+ = \mathbf{L}_{\text{pop}}^+ - \eta \mathbf{E}_{\mathbf{L}}, \quad \mathbf{R}^+ = \mathbf{R}_{\text{pop}}^+ - \eta \mathbf{E}_{\mathbf{R}}. \quad (63)$$

Here,  $(\mathbf{L}_{\text{pop}}^+, \mathbf{R}_{\text{pop}}^+)$  represents the population-level update from  $\Omega_1^*$

$$\mathbf{L}_{\text{pop}}^+ := \mathbf{L} - \eta p_1 w_1 \Delta_1 \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1}, \quad \mathbf{R}_{\text{pop}}^+ := \mathbf{R} - \eta p_1 w_1 \Delta_1^\top \mathbf{L} (\mathbf{L}^\top \mathbf{L})^{-1}, \quad (64)$$

and the residual components are given by

$$\mathbf{E}_{\mathbf{L}} := (\Delta_{\text{mis}} + \Delta_{\text{fs}}) \mathbf{R} (\mathbf{R}^\top \mathbf{R})^{-1}, \quad \mathbf{E}_{\mathbf{R}} := (\Delta_{\text{mis}} + \Delta_{\text{fs}})^\top \mathbf{L} (\mathbf{L}^\top \mathbf{L})^{-1}$$

with

$$\Delta_{\text{mis}} := \sum_{k \neq 1} p_k w_k \Delta_k, \quad \Delta_{\text{fs}} := \sum_{k=1}^K p_k \left( \frac{1}{|\Omega_k^*|} \sum_{i \in \Omega_k^*} (\langle \mathbf{A}_i, \Delta_k \rangle - \zeta_i) \mathbb{1}(\langle \mathbf{A}_i, \Delta_k \rangle - \zeta_i; \tau) \mathbf{A}_i - w_k \Delta_k \right). \quad (65)$$

Before moving on, we note that it is crucial to control the sizes of  $\Delta_{\text{mis}}$  and  $\Delta_{\text{fs}}$ , where “mis” stands for “mislabeling”, and “fs” stands for “finite sample”. Regarding  $\Delta_{\text{mis}}$ , Fact 1 tells us that for all  $k \neq 1$ ,

$$w_k \|\Delta_k\|_{\text{F}} = w_1 \frac{w_k}{w_1} \|\Delta_k\|_{\text{F}} \leq w_1 \frac{\|\Delta_1\|_{\text{F}}^2 + \sigma^2}{\|\Delta_k\|_{\text{F}}^2 + \sigma^2} \|\Delta_k\|_{\text{F}} \leq \frac{\|\Delta_1\|_{\text{F}}^2 + \sigma^2}{\|\Delta_k\|_{\text{F}}} \leq 2 \frac{\|\Delta_1\|_{\text{F}}^2 + \sigma^2}{\|\mathbf{M}_k^* - \mathbf{M}_1^*\|_{\text{F}}}.$$

Here, the last inequality holds since

$$\|\Delta_k\|_{\text{F}} = \|\mathbf{L}\mathbf{R}^\top - \mathbf{M}_k^*\|_{\text{F}} \geq \|\mathbf{M}_1^* - \mathbf{M}_k^*\|_{\text{F}} - \|\Delta_1\|_{\text{F}} \geq 0.5\|\mathbf{M}_1^* - \mathbf{M}_k^*\|_{\text{F}},$$

where we have used  $\|\Delta_1\|_{\text{F}} \leq c_2 \sigma_{r_1}(\mathbf{M}_1^*) \leq 0.5\|\mathbf{M}_1^* - \mathbf{M}_k^*\|_{\text{F}}$  due to the assumption that  $(\mathbf{L}, \mathbf{R}) \in \mathcal{B}$  defined in (56). Consequently, we obtain

$$\|\Delta_{\text{mis}}\|_{\text{F}} = \left\| \sum_{k \neq 1} p_k w_k \Delta_k \right\|_{\text{F}} \leq \sum_{k \neq 1} p_k w_k \|\Delta_k\|_{\text{F}} \leq 2 \frac{\|\Delta_1\|_{\text{F}}^2 + \sigma^2}{\min_{k \neq 1} \|\mathbf{M}_k^* - \mathbf{M}_1^*\|_{\text{F}}}. \quad (66)$$

Next, we turn to the term  $\Delta_{\text{fs}}$ . Note that  $\text{rank}(\Delta_k) \leq 2r$ . Therefore, Lemmas 1 and 2 (see Remark 2) imply that, with probability at least  $1 - Ce^{-cn}$  for some constants  $c, C > 0$ , the following holds simultaneously for all  $(\mathbf{L}, \mathbf{R}) \in \mathcal{B}$  (cf. (56)):

1. the truncating level  $\tau$  obeys

$$0.54 < \frac{\tau}{\sqrt{\|\Delta_1\|_{\text{F}}^2 + \sigma^2}} < 1.35; \quad (67)$$

2. for any real matrix  $\mathbf{W}$  with  $n_2$  rows and of rank at most  $r$ , we have

$$\|\Delta_{\text{fs}} \mathbf{W}\|_{\text{F}} \leq \sum_{k=1}^K p_k \delta \tau \|\mathbf{W}\| = \delta \tau \|\mathbf{W}\| \leq 1.35 \delta \sqrt{\|\Delta_1\|_{\text{F}}^2 + \sigma^2} \|\mathbf{W}\|. \quad (68)$$

The above-mentioned bounds will play a useful role in subsequent steps.

**Step 2: per-iteration improvement above the error floor (34).** Let us look at the Euclidean error

$$\|\mathbf{L}^+(\mathbf{R}^+)^\top - \mathbf{M}_1^*\|_{\text{F}} = \|(\mathbf{L}_{\text{pop}}^+ - \eta \mathbf{E}_L)(\mathbf{R}_{\text{pop}}^+ - \eta \mathbf{E}_R)^\top - \mathbf{M}_1^*\|_{\text{F}} \quad (69a)$$

$$\leq \|\mathbf{L}_{\text{pop}}^+(\mathbf{R}_{\text{pop}}^+)^\top - \mathbf{M}_1^*\|_{\text{F}} + \eta \left( \|\mathbf{E}_L(\mathbf{R}_{\text{pop}}^+)^\top\|_{\text{F}} + \|\mathbf{L}_{\text{pop}}^+(\mathbf{E}_R)^\top\|_{\text{F}} + \eta \|\mathbf{E}_L(\mathbf{E}_R)^\top\|_{\text{F}} \right). \quad (69b)$$

Since  $\mathbf{L}_{\text{pop}}^+$  and  $\mathbf{R}_{\text{pop}}^+$  (64) are exactly the same as the update rule of scaled gradient descent for low-rank matrix factorization, [TMC20, Theorem 3] tells us that if  $0 < \eta p_1 w_1 \leq 2/3$  (which holds true under our choices of  $\eta \leq 1.3/p_1$  and  $\alpha \leq 0.8p_1$ ), then

$$\|\mathbf{L}_{\text{pop}}^+(\mathbf{R}_{\text{pop}}^+)^\top - \mathbf{M}_1^*\|_{\text{F}} \leq (1 - 0.7\eta p_1 w_1) \|\mathbf{L}\mathbf{R}^\top - \mathbf{M}_1^*\|_{\text{F}}. \quad (70)$$

It remains to control the perturbation terms in (69b), accomplished as follows.

**Claim 4.** Denoting

$$B := 2 \left( \delta \sqrt{\|\Delta_1\|_{\text{F}}^2 + \sigma^2} + \frac{\|\Delta_1\|_{\text{F}}^2 + \sigma^2}{\min_{k \neq 1} \|\mathbf{M}_k^* - \mathbf{M}_1^*\|_{\text{F}}} \right), \quad (71)$$

one has

$$\max \left\{ \|\mathbf{E}_L(\mathbf{R}_{\text{pop}}^+)^\top\|_{\text{F}}, \|\mathbf{L}_{\text{pop}}^+(\mathbf{E}_R)^\top\|_{\text{F}} \right\} \leq 2B, \quad \|\mathbf{E}_L(\mathbf{E}_R)^\top\|_{\text{F}} \leq \frac{2}{\sigma_{r_1}(\mathbf{M}_1^*)} B^2. \quad (72)$$

Putting (70) and (72) back to (69) and denoting  $\Delta_1^+ := \mathbf{L}^+(\mathbf{R}^+)^\top - \mathbf{M}_1^*$ , we have

$$\|\Delta_1^+\|_{\text{F}} \leq (1 - 0.7\eta p_1 w_1) \|\Delta_1\|_{\text{F}} + \eta \left( 4B + \frac{2\eta}{\sigma_{r_1}(\mathbf{M}_1^*)} B^2 \right). \quad (73)$$

It remains to control  $B$ . First, the relations  $\delta \leq c_0/K$  and  $\|\Delta_1\|_F \geq C_2 K \delta \sigma$  (for some sufficiently large constant  $C_2 > 0$ ) imply that

$$\delta \sqrt{\|\Delta_1\|_F + \sigma^2} \leq \delta \|\Delta_1\|_F + \delta \sigma \leq \frac{c_3}{K} \|\Delta_1\|_F \quad (74)$$

for some sufficiently small constant  $c_3 > 0$ . Moreover, observing that

$$\frac{C_2 K \sigma^2}{\min_{k \neq 1} \|\mathbf{M}_k^* - \mathbf{M}_1^*\|_F} \leq \|\Delta_1\|_F \leq \frac{c_1}{K} \min_{k \neq 1} \|\mathbf{M}_k^* - \mathbf{M}_1^*\|_F,$$

we have

$$\frac{\|\Delta_1\|_F^2 + \sigma^2}{\min_{k \neq 1} \|\mathbf{M}_k^* - \mathbf{M}_1^*\|_F} = \frac{\|\Delta_1\|_F^2}{\min_{k \neq 1} \|\mathbf{M}_k^* - \mathbf{M}_1^*\|_F} + \frac{\sigma^2}{\min_{k \neq 1} \|\mathbf{M}_k^* - \mathbf{M}_1^*\|_F} \leq \frac{c_4}{K} \|\Delta_1\|_F \quad (75)$$

for some sufficiently small constant  $c_4 > 0$ . Putting (74) and (75) back into (71), we have

$$B \leq \frac{2(c_3 + c_4)}{K} \|\Delta_1\|_F, \quad (76)$$

which together with  $\|\Delta_1\|_F \leq c_1 \sigma_{r_1}(\mathbf{M}_1^*)$  implies the existence of some small constant  $c_5 > 0$  such that

$$4B + \frac{2\eta}{\sigma_{r_1}(\mathbf{M}_1^*)} B^2 = 4B \left( 1 + \frac{\eta}{2\sigma_{r_1}(\mathbf{M}_1^*)} B \right) \leq 8B \leq \frac{c_5}{K} \|\Delta_1\|_F.$$

Substituting this into (73), we arrive at the desired bound

$$\|\Delta_1^+\|_F \leq (1 - c_2 \eta p_1) \|\Delta_1\|_F$$

for some constant  $c_2 > 0$ ; this is because in (73), we have  $p_1 \asymp 1/K$  by assumption, and  $w_1 \gtrsim 1$  according to (67).

**Step 3: no blowing up below the error floor (34).** Suppose that the estimation error satisfies

$$\|\Delta_1\|_F \lesssim \max \left\{ K \sigma \delta, \frac{K \sigma^2}{\min_{k \neq 1} \|\mathbf{M}_k^* - \mathbf{M}_1^*\|_F} \right\}. \quad (77)$$

We intend to show that, in this case, the estimation error of the next iterate  $\|\Delta_1^+\|_F$  satisfies the same upper bound (77); if this claim were true, then combining this with our results in Step 2 would complete the convergence analysis of ScaledTGD.

Note that (73) remains valid when  $\|\Delta_1\|_F$  is below the error floor, which implies that

$$\|\Delta_1^+\|_F \lesssim \|\Delta_1\|_F + KB \left( 1 + \frac{KB}{\sigma_{r_1}(\mathbf{M}_1^*)} \right). \quad (78)$$

Recalling the definition of  $B$  in (71), one has

$$KB \lesssim K(\|\Delta_1\|_F + \sigma) \left( \delta + \frac{\|\Delta_1\|_F + \sigma}{\min_{k \neq 1} \|\mathbf{M}_k^* - \mathbf{M}_1^*\|_F} \right).$$

By the assumption that  $\delta \lesssim 1/K$  and  $\sigma \lesssim \min_k \|\mathbf{M}_k^*\|_F / K$ , we have  $\|\Delta_1\|_F \lesssim \sigma$  according to (77), and thus  $KB / \sigma_{r_1}(\mathbf{M}_1^*) \lesssim \sigma / \sigma_{r_1}(\mathbf{M}_1^*) \lesssim 1$ . Consequently, on the right-hand side of (78) we have

$$KB \left( 1 + \frac{KB}{\sigma_{r_1}(\mathbf{M}_1^*)} \right) \lesssim KB \lesssim K \sigma \left( \delta + \frac{\sigma}{\min_{k \neq 1} \|\mathbf{M}_k^* - \mathbf{M}_1^*\|_F} \right),$$

which has exactly the same form as the error floor in (77). This completes our proof for this step.

*Proof of Claim 4.* We shall only prove the first part of (72) concerning  $\|\mathbf{E}_L(\mathbf{R}_{\text{pop}}^+)^{\top}\|_{\text{F}}$ ; the analysis for  $\|\mathbf{L}_{\text{pop}}^+(\mathbf{E}_R)^{\top}\|_{\text{F}}$  is essentially the same. By the triangle inequality, we have

$$\|\mathbf{E}_L(\mathbf{R}_{\text{pop}}^+)^{\top}\|_{\text{F}} \leq \|\mathbf{E}_L\mathbf{R}^{\top}\|_{\text{F}} + \eta p_1 w_1 \|\mathbf{E}_L(\mathbf{L}^{\top}\mathbf{L})^{-1}\mathbf{L}^{\top}\mathbf{\Delta}_1\|_{\text{F}}. \quad (79)$$

We utilize (66) and (68) from Step 1 to control the terms above. For the first term of (79), recognizing that  $\|\mathbf{R}(\mathbf{R}^{\top}\mathbf{R})^{-1}\mathbf{R}^{\top}\| \leq 1$  we have

$$\begin{aligned} \|\mathbf{E}_L\mathbf{R}^{\top}\|_{\text{F}} &= \|(\mathbf{\Delta}_{\text{fs}} + \mathbf{\Delta}_{\text{mis}})\mathbf{R}(\mathbf{R}^{\top}\mathbf{R})^{-1}\mathbf{R}^{\top}\|_{\text{F}} \leq \|\mathbf{\Delta}_{\text{fs}}\mathbf{R}(\mathbf{R}^{\top}\mathbf{R})^{-1}\mathbf{R}^{\top}\|_{\text{F}} + \|\mathbf{\Delta}_{\text{mis}}\|_{\text{F}} \\ &\leq 2 \left( \delta \sqrt{\|\mathbf{\Delta}_1\|_{\text{F}}^2 + \sigma^2} + \frac{\|\mathbf{\Delta}_1\|_{\text{F}}^2 + \sigma^2}{\min_{k \neq 1} \|\mathbf{M}_k^* - \mathbf{M}_1^*\|_{\text{F}}} \right) = B. \end{aligned}$$

Regarding the second term of (79), we observe that

$$\begin{aligned} \|\mathbf{E}_L(\mathbf{L}^{\top}\mathbf{L})^{-1}\mathbf{L}^{\top}\mathbf{\Delta}_1\|_{\text{F}} &= \|(\mathbf{\Delta}_{\text{fs}} + \mathbf{\Delta}_{\text{mis}})\mathbf{R}(\mathbf{R}^{\top}\mathbf{R})^{-1}(\mathbf{L}^{\top}\mathbf{L})^{-1}\mathbf{L}^{\top}\mathbf{\Delta}_1\|_{\text{F}} \\ &\leq \left( \|\mathbf{\Delta}_{\text{fs}}\mathbf{R}(\mathbf{R}^{\top}\mathbf{R})^{-1}(\mathbf{L}^{\top}\mathbf{L})^{-1}\mathbf{L}^{\top}\|_{\text{F}} + \|\mathbf{\Delta}_{\text{mis}}\|_{\text{F}} \|\mathbf{R}(\mathbf{R}^{\top}\mathbf{R})^{-1}(\mathbf{L}^{\top}\mathbf{L})^{-1}\mathbf{L}^{\top}\| \right) \|\mathbf{\Delta}_1\|_{\text{F}} \\ &\stackrel{(i)}{\leq} 2 \left( \delta \sqrt{\|\mathbf{\Delta}_1\|_{\text{F}}^2 + \sigma^2} + \frac{\|\mathbf{\Delta}_1\|_{\text{F}}^2 + \sigma^2}{\min_{k \neq 1} \|\mathbf{M}_k^* - \mathbf{M}_1^*\|_{\text{F}}} \right) \frac{2}{\sigma_{r_1}(\mathbf{M}_1^*)} \cdot c_1 \sigma_{r_1}(\mathbf{M}_1^*) \\ &= 4c_1 \left( \delta \sqrt{\|\mathbf{\Delta}_1\|_{\text{F}}^2 + \sigma^2} + \frac{\|\mathbf{\Delta}_1\|_{\text{F}}^2 + \sigma^2}{\min_{k \neq 1} \|\mathbf{M}_k^* - \mathbf{M}_1^*\|_{\text{F}}} \right) = 2c_1 B, \end{aligned}$$

where (i) follows from  $\|\mathbf{\Delta}_1\|_{\text{F}} \leq c_1 \sigma_{r_1}(\mathbf{M}_1^*)$  (see (56)) as well as the following fact (which will be proved at the end of this section): for any  $\mathbf{L} \in \mathbb{R}^{n_1 \times r_1}$  and  $\mathbf{R} \in \mathbb{R}^{n_2 \times r_1}$ ,

$$\text{if } \|\mathbf{L}\mathbf{R}^{\top} - \mathbf{M}_1^*\|_{\text{F}} \leq \frac{\sigma_{r_1}(\mathbf{M}_1^*)}{2}, \text{ then } \|\mathbf{L}(\mathbf{L}^{\top}\mathbf{L})^{-1}(\mathbf{R}^{\top}\mathbf{R})^{-1}\mathbf{R}^{\top}\| \leq \frac{2}{\sigma_{r_1}(\mathbf{M}_1^*)}. \quad (80)$$

Combining these with (79) establishes that  $\|\mathbf{E}_L(\mathbf{R}_{\text{pop}}^+)^{\top}\|_{\text{F}} \leq 2B$ , which is the first part of (72).

Finally, for the second part of (72), we can apply similar techniques to reach

$$\begin{aligned} \|\mathbf{E}_L(\mathbf{E}_R)^{\top}\|_{\text{F}} &= \|(\mathbf{\Delta}_{\text{fs}} + \mathbf{\Delta}_{\text{mis}})\mathbf{R}(\mathbf{R}^{\top}\mathbf{R})^{-1}(\mathbf{L}^{\top}\mathbf{L})^{-1}\mathbf{L}^{\top}(\mathbf{\Delta}_{\text{fs}} + \mathbf{\Delta}_{\text{mis}})\|_{\text{F}} \\ &\leq \frac{2}{\sigma_{r_1}(\mathbf{M}_1^*)} \cdot 4 \left( \delta \sqrt{\|\mathbf{\Delta}_1\|_{\text{F}}^2 + \sigma^2} + \frac{\|\mathbf{\Delta}_1\|_{\text{F}}^2 + \sigma^2}{\min_{k \neq 1} \|\mathbf{M}_k^* - \mathbf{M}_1^*\|_{\text{F}}} \right)^2 = \frac{2}{\sigma_{r_1}(\mathbf{M}_1^*)} B^2. \end{aligned}$$

□

*Proof of (80).* Weyl's inequality tells us that

$$\sigma_{r_1}(\mathbf{L}\mathbf{R}^{\top}) \geq \sigma_{r_1}(\mathbf{M}_1^*) - \|\mathbf{L}\mathbf{R}^{\top} - \mathbf{M}_1^*\|_{\text{F}} \geq \frac{\sigma_{r_1}(\mathbf{M}_1^*)}{2}, \quad (81)$$

which further implies that both  $\mathbf{L}$  and  $\mathbf{R}$  have full column rank  $r_1$ . Consequently, we denote the SVD of  $\mathbf{L}$  and  $\mathbf{R}$  as  $\mathbf{L} = \mathbf{U}_L \mathbf{\Sigma}_L \mathbf{V}_L^{\top}$  and  $\mathbf{R} = \mathbf{U}_R \mathbf{\Sigma}_R \mathbf{V}_R^{\top}$ , where  $\mathbf{V}_L, \mathbf{V}_R$  are  $r_1 \times r_1$  orthonormal matrices. With the SVD representations in place, it is easy to check that

$$\mathbf{L}\mathbf{R}^{\top} = \mathbf{U}_L \mathbf{\Sigma}_L \mathbf{V}_L^{\top} \mathbf{V}_R \mathbf{\Sigma}_R \mathbf{U}_R^{\top}, \quad \text{and} \quad \mathbf{L}(\mathbf{L}^{\top}\mathbf{L})^{-1}(\mathbf{R}^{\top}\mathbf{R})^{-1}\mathbf{R}^{\top} = \mathbf{U}_L \mathbf{\Sigma}_L^{-1} \mathbf{V}_L^{\top} \mathbf{V}_R \mathbf{\Sigma}_R^{-1} \mathbf{U}_R^{\top}.$$

In addition, the orthonormality of  $\mathbf{V}_L$  and  $\mathbf{V}_R$  implies

$$(\mathbf{\Sigma}_L \mathbf{V}_L^{\top} \mathbf{V}_R \mathbf{\Sigma}_R)^{-1} = \mathbf{\Sigma}_R^{-1} (\mathbf{V}_L^{\top} \mathbf{V}_R)^{-1} \mathbf{\Sigma}_L^{-1} = \mathbf{\Sigma}_R^{-1} \mathbf{V}_R^{\top} \mathbf{V}_L \mathbf{\Sigma}_L^{-1} = (\mathbf{\Sigma}_L^{-1} \mathbf{V}_L^{\top} \mathbf{V}_R \mathbf{\Sigma}_R^{-1})^{\top},$$

thus indicating that

$$\begin{aligned} \|\mathbf{L}(\mathbf{L}^{\top}\mathbf{L})^{-1}(\mathbf{R}^{\top}\mathbf{R})^{-1}\mathbf{R}^{\top}\| &= \|\mathbf{U}_L \mathbf{\Sigma}_L^{-1} \mathbf{V}_L^{\top} \mathbf{V}_R \mathbf{\Sigma}_R^{-1} \mathbf{U}_R^{\top}\| = \|\mathbf{\Sigma}_L^{-1} \mathbf{V}_L^{\top} \mathbf{V}_R \mathbf{\Sigma}_R^{-1}\| \\ &= \|(\mathbf{\Sigma}_L \mathbf{V}_L^{\top} \mathbf{V}_R \mathbf{\Sigma}_R)^{-1}\| = \frac{1}{\sigma_{r_1}(\mathbf{\Sigma}_L \mathbf{V}_L^{\top} \mathbf{V}_R \mathbf{\Sigma}_R)} = \frac{1}{\sigma_{r_1}(\mathbf{L}\mathbf{R}^{\top})}. \end{aligned}$$

Combining this with (81) completes the proof. □

## C Technical lemmas

This section collects several technical lemmas that are helpful for our analysis (particularly for the analysis of Stage 3). For notational convenience, we define the set of low-rank matrices as

$$\mathcal{R}_r := \{\mathbf{X} \in \mathbb{R}^{n_1 \times n_2} : \text{rank}(\mathbf{X}) \leq r\}. \quad (82)$$

We remind the reader of the definitions  $\mathbb{1}(a; b) = \mathbb{1}(|a| \leq b)$  for  $a, b \in \mathbb{R}$  and  $w(x) = \int_{-x}^x t^2 \phi(t) dt$  for  $x \geq 0$ .

**Variants of matrix-RIP.** We recall the standard notion of restricted isometry property (RIP) from the literature of matrix sensing, and introduce a variant called *truncated* RIP (TRIP).

**Definition 2.** Let  $\{\mathbf{A}_i\}_{i=1}^m$  be a set of matrices in  $\mathbb{R}^{n_1 \times n_2}$ . Consider  $1 \leq r \leq \min\{n_1, n_2\}$  and  $0 < \delta < 1$ .

1. We say that  $\{\mathbf{A}_i\}_{1 \leq i \leq m}$  satisfy  $(r, \delta)$ -RIP if

$$\left| \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle - \langle \mathbf{X}, \mathbf{Z} \rangle \right| \leq \delta \|\mathbf{X}\|_{\text{F}} \|\mathbf{Z}\|_{\text{F}} \quad (83)$$

holds simultaneously for all  $\mathbf{X}, \mathbf{Z} \in \mathcal{R}_r$ .

2. We say that  $\{\mathbf{A}_i\}_{1 \leq i \leq m}$  satisfy  $(r, \delta)$ -TRIP if

$$\left| \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X} \rangle \mathbb{1}(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau \|\mathbf{X}\|_{\text{F}}) \langle \mathbf{A}_i, \mathbf{Z} \rangle - w(\tau) \langle \mathbf{X}, \mathbf{Z} \rangle \right| \leq \delta \tau \|\mathbf{X}\|_{\text{F}} \|\mathbf{Z}\|_{\text{F}} \quad (84)$$

holds simultaneously for all  $\mathbf{X}, \mathbf{Z} \in \mathcal{R}_r$  and for all  $0 \leq \tau \leq 1.35$ .

As it turns out, the Gaussian design satisfies the above notion of RIP and TRIP, as formalized below.

**Lemma 1.** Let  $\{\mathbf{A}_i\}_{1 \leq i \leq m}$  be random matrices in  $\mathbb{R}^{n_1 \times n_2}$  with i.i.d.  $\mathcal{N}(0, 1)$  entries, and denote  $n := \max\{n_1, n_2\}$ . There exist some sufficiently large constants  $C_1, C_3 > 0$  and some other constants  $C_2, c_2, C_4, c_4 > 0$  such that

1. If  $m \geq C_1 n r \delta^{-2} \log(1/\delta)$ , then with probability at least  $1 - C_2 e^{-c_2 n}$ ,  $\{\mathbf{A}_i\}_{1 \leq i \leq m}$  satisfy  $(r, \delta)$ -RIP.
2. If  $m \geq C_3 n r \delta^{-2} \log m$ , then with probability at least  $1 - C_4 e^{-c_4 n}$ ,  $\{\mathbf{A}_i\}_{1 \leq i \leq m}$  satisfy  $(r, \delta)$ -TRIP.

**Empirical quantiles.** Our next technical lemma is a uniform concentration result for empirical quantiles. Given the design matrices  $\{\mathbf{A}_i\}_{1 \leq i \leq N}$ , the index sets  $\{\Omega_k^*\}_{1 \leq k \leq K}$  and the low-rank matrices  $\{\mathbf{X}_k\}_{1 \leq k \leq K}$ , we define several sets as follows:

$$\mathcal{D}_k := \left\{ |\langle \mathbf{A}_i, \mathbf{X}_k \rangle| \right\}_{i \in \Omega_k^*}, \quad 1 \leq k \leq K; \quad \mathcal{D} := \mathcal{D}_1 \cup \dots \cup \mathcal{D}_K. \quad (85)$$

In addition, let us introduce the following set of low-rank matrices:

$$\mathcal{T}_1 := \left\{ (\mathbf{X}_1, \dots, \mathbf{X}_K) : \mathbf{X}_k \in \mathcal{R}_r, 1 \leq k \leq K; 0 < \|\mathbf{X}_1\|_{\text{F}} \leq \frac{c_0}{K} \min_{k \neq 1} \|\mathbf{X}_k\|_{\text{F}} \right\}, \quad (86)$$

where  $c_0 > 0$  is some sufficiently small constant. Recall that  $Q_\alpha(\mathcal{D})$  denotes the  $\alpha$ -quantile of  $\mathcal{D}$ , as defined in (2).

**Lemma 2.** Let  $\{\mathbf{A}_i\}_{1 \leq i \leq N}$  be random matrices in  $\mathbb{R}^{n_1 \times n_2}$  with i.i.d.  $\mathcal{N}(0, 1)$  entries. Set  $n = \max\{n_1, n_2\}$ , and suppose the index sets  $\{\Omega_k^*\}_{1 \leq k \leq K}$  are disjoint and satisfy the condition (18). If  $0.6p_1 \leq \alpha \leq 0.8p_1$  and  $N \geq C_0 n r K^3 \log N$  for some sufficiently large constant  $C_0 > 0$ , then there exist some universal constants  $C, c > 0$  such that: with probability at least  $1 - C e^{-cn}$ ,

$$0.54 < \frac{Q_\alpha(\mathcal{D})}{\|\mathbf{X}_1\|_{\text{F}}} < 1.35$$

holds simultaneously for all  $(\mathbf{X}_1, \dots, \mathbf{X}_K) \in \mathcal{T}_1$ , where  $\mathcal{D}$  is defined in (85).



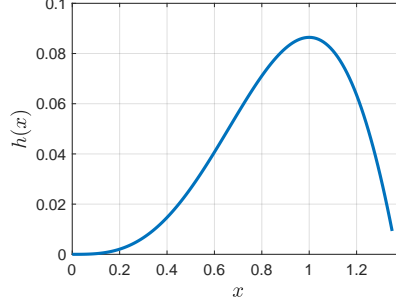


Figure 2: The function  $h(\cdot)$  defined in (89) is nonnegative over the interval  $(0, 1.35]$ .

*Remark 2.* We can further incorporate additional Gaussian noise  $\{\zeta_i\}$  into Lemmas 1 and 2, where  $\zeta_i$  i.i.d.  $\mathcal{N}(0, \sigma^2)$ . For example, we claim that, with the same sample complexity  $m$  as in Lemma 1, we have the following noisy version of  $(r, \delta)$ -TRIP (84):

$$\left| \frac{1}{m} \sum_{i=1}^m (\langle \mathbf{A}_i, \mathbf{X} \rangle - \zeta_i) \mathbb{1} \left( \langle \mathbf{A}_i, \mathbf{X} \rangle - \zeta_i; \tau \sqrt{\|\mathbf{X}\|_{\mathbb{F}}^2 + \sigma^2} \right) \langle \mathbf{A}_i, \mathbf{Z} \rangle - w(\tau) \langle \mathbf{X}, \mathbf{Z} \rangle \right| \leq \delta \tau \sqrt{\|\mathbf{X}\|_{\mathbb{F}}^2 + \sigma^2} \|\mathbf{Z}\|_{\mathbb{F}}. \quad (87)$$

To see this, let us define the augmented matrices

$$\mathbf{X}^{\text{aug}} := \begin{bmatrix} \mathbf{X} & \mathbf{0} \\ \mathbf{0} & -\sigma \end{bmatrix}, \quad \mathbf{Z}^{\text{aug}} := \begin{bmatrix} \mathbf{Z} & \mathbf{0} \\ \mathbf{0} & 0 \end{bmatrix}, \quad \mathbf{A}_i^{\text{aug}} := \begin{bmatrix} \mathbf{A}_i & * \\ * & \zeta_i/\sigma \end{bmatrix}, \quad 1 \leq i \leq m,$$

where  $*$  stands for some auxiliary i.i.d.  $\mathcal{N}(0, 1)$  entries. Observe that  $\{\mathbf{A}_i^{\text{aug}}\}_{1 \leq i \leq m}$  are random matrices with i.i.d.  $\mathcal{N}(0, 1)$  entries; in addition,  $\text{rank}(\mathbf{X}^{\text{aug}}) = \text{rank}(\mathbf{X}) + 1$ ,  $\text{rank}(\mathbf{Z}^{\text{aug}}) = \text{rank}(\mathbf{Z})$ , and  $\|\mathbf{X}^{\text{aug}}\|_{\mathbb{F}}^2 = \|\mathbf{X}\|_{\mathbb{F}}^2 + \sigma^2$ ; finally,  $\langle \mathbf{A}_i, \mathbf{X} \rangle - \zeta_i = \langle \mathbf{A}_i^{\text{aug}}, \mathbf{X}^{\text{aug}} \rangle$ ,  $\langle \mathbf{A}_i, \mathbf{Z} \rangle = \langle \mathbf{A}_i^{\text{aug}}, \mathbf{Z}^{\text{aug}} \rangle$ , and  $\langle \mathbf{X}, \mathbf{Z} \rangle = \langle \mathbf{X}^{\text{aug}}, \mathbf{Z}^{\text{aug}} \rangle$ . Therefore, the left-hand side of (87) can be equivalently written as in the noiseless form (84), in terms of these augmented matrices, thus allowing us to apply Lemma 1 to prove (87). This trick of augmentation can be applied to Lemma 2 as well, which we omit here for brevity.

**One miscellaneous result.** Further, we record below a basic property concerning the function  $w(\cdot)$ .

**Fact 1.** *The function  $w(\cdot)$  defined in (60) satisfies*

$$\frac{w(x)}{w(y)} \leq \frac{x^2}{y^2}, \quad 0 < x \leq y \leq 1.35. \quad (88)$$

*Proof.* This result is equivalent to saying  $w(x)/x^2 \leq w(y)/y^2$  for any  $0 < x \leq y \leq 1.35$ . Hence, it suffices to show that the function  $g(x) := w(x)/x^2$  is nondecreasing over  $(0, 1.35]$ , or equivalently,

$$h(x) := \sqrt{\frac{2}{\pi}} x^3 e^{-\frac{x^2}{2}} - 2w(x), \quad g'(x) = \frac{1}{x^3} h(x) \geq 0, \quad 0 < x \leq 1.35. \quad (89)$$

This can be verified numerically (see Figure 2), which completes the proof.  $\square$

The rest of this section is devoted to proving Lemmas 1 and 2. We use the standard notions (e.g. the subgaussian norm  $\|\cdot\|_{\psi_2}$ ) and properties related to subgaussian random variables (cf. [Ver18, Section 2]). For notational convenience, we define the normalized version of  $\mathcal{R}_r$  defined in (82), as follows:

$$\mathcal{R}_r^{\text{norm}} := \{ \mathbf{X} \in \mathbb{R}^{n_1 \times n_2} : \text{rank}(\mathbf{X}) \leq r, \|\mathbf{X}\|_{\mathbb{F}} = 1 \}. \quad (90)$$

Before moving on, we record two results that will be useful throughout the proof.

**Lemma 3.** Let  $\{\mathbf{A}_i\}_{i=1}^m$  be a set of random matrices in  $\mathbb{R}^{n_1 \times n_2}$  with i.i.d.  $\mathcal{N}(0, 1)$  entries. Denote  $n := \max\{n_1, n_2\}$ , and let  $Z$  be a random variable having the same distribution as  $|\mathcal{N}(0, 1)|$ . For all  $t > 0$  and  $0 < \epsilon < 1$ , with probability at least  $1 - (9/\epsilon)^{3nr} \exp(-c_1 mt^2/(\tau + t)) - C_2 e^{-c_2 n}$ , the following

$$\frac{1}{m} \sum_{i=1}^m \mathbb{1}(|\langle \mathbf{A}_i, \mathbf{X} \rangle| \leq \tau) \leq \mathbb{P}(Z \leq 1.01\tau) + t + \frac{200\epsilon}{\tau}$$

holds simultaneously for all  $\mathbf{X} \in \mathcal{R}_r^{\text{norm}}$ , provided that  $m \geq Cnr \log m$ . Here,  $c_1, C_2, c_2 > 0$  are universal constants, and  $C > 0$  is some sufficiently large constant.

**Proposition 2.** Consider  $\mathbf{a}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ ,  $1 \leq i \leq m$ . There exist some universal constants  $C, c > 0$  such that with probability at least  $1 - Ce^{-cd}$ , we have

$$\max_{1 \leq i \leq m} \|\mathbf{a}_i\|_2 \lesssim \sqrt{d} + \sqrt{\log m}.$$

*Proof.* This result follows from [Ver18, Corollary 7.3.3] and the union bound.  $\square$

## C.1 Proof of Lemma 1

The first result on RIP has been established in the literature (e.g. [CP11, Theorem 2.3]), and hence we only need to prove the second result on TRIP. We first restrict to the case

$$m^{-100} \leq \tau \leq 1.35;$$

at the end of this subsection, we will prove TRIP for the case  $0 \leq \tau < m^{-100}$  separately. By homogeneity, it is sufficient to show that

$$\left| \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X} \rangle \mathbb{1}(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau) \langle \mathbf{A}_i, \mathbf{Z} \rangle - w(\tau) \langle \mathbf{X}, \mathbf{Z} \rangle \right| \leq \delta \tau \quad (91)$$

holds simultaneously for all  $(\mathbf{X}, \mathbf{Z}, \tau) \in \mathcal{T}_{\text{TRIP}}$ , where

$$\mathcal{T}_{\text{TRIP}} := \{(\mathbf{X}, \mathbf{Z}, \tau) : \mathbf{X}, \mathbf{Z} \in \mathcal{R}_r^{\text{norm}}, m^{-100} \leq \tau \leq 1.35\}.$$

The proof consists of two steps: (1) we replace the discontinuous function  $\mathbb{1}$  by a Lipschitz continuous surrogate  $\chi$  and establish a uniform concentration result for  $\chi$ ; (2) we show that the discrepancy incurred by replacing  $\mathbb{1}$  with  $\chi$  is uniformly small. Our proof argument is conditioned on the high-probability event that  $\{\mathbf{A}_i\}_{i=1}^m$  satisfy  $(2r, \delta)$ -RIP.

**Step 1: replacing  $\mathbb{1}$  with  $\chi$ .** Define an auxiliary function  $\chi$  as follows: for all  $a \in \mathbb{R}$  and  $\tau > 0$ ,

$$\chi(a; \tau) := \begin{cases} 1, & |a| \leq (1 - c_\chi)\tau; \\ 0, & |a| \geq \tau; \\ \frac{\tau - |a|}{c_\chi \tau}, & (1 - c_\chi)\tau < |a| < \tau. \end{cases} \quad (92)$$

Here we set the parameter

$$c_\chi := c_0 \delta^2 m^{-100} \quad (93)$$

for some sufficiently small constant  $c_0 > 0$ , whose rationale will be made apparent in Step 2. It is easily seen that  $\chi$  enjoys the following properties:

- (Continuity) For any  $\tau > 0$ ,  $\chi(\cdot; \tau)$  is piecewise linear and  $1/(c_\chi \tau)$ -Lipschitz continuous.
- (Closeness to  $\mathbb{1}$ ) For any  $\tau > 0$  and  $a \in \mathbb{R}$ ,  $\chi(a; \tau) \leq \mathbb{1}(a; \tau) \leq \chi(a; \tau/(1 - c_\chi))$ .
- (Homogeneity) For any  $\tau > 0$ ,  $a \in \mathbb{R}$  and  $c_0 > 0$ ,  $\chi(a; \tau) = \chi(a/c_0; \tau/c_0)$ .

- If  $0 \leq \epsilon_\tau \leq c_\chi \tau$  and  $\tau - \epsilon_\tau \leq \tau_0 \leq \tau$ , then  $\|\chi(\cdot; \tau) - \chi(\cdot; \tau_0)\|_\infty = \chi(\tau_0; \tau) = (\tau - \tau_0)/(c_\chi \tau) \leq \epsilon_\tau/(c_\chi \tau)$ .
- The function  $f(a) := a \cdot \chi(a; \tau)$  is  $1/c_\chi$ -Lipschitz continuous.

For notational convenience, define

$$E_m^\chi(\mathbf{X}, \mathbf{Z}, \tau) := \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X} \rangle \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau) \langle \mathbf{A}_i, \mathbf{Z} \rangle, \quad (94a)$$

$$E^\chi(\mathbf{X}, \mathbf{Z}, \tau) := \mathbb{E}[\langle \mathbf{A}_i, \mathbf{X} \rangle \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau) \langle \mathbf{A}_i, \mathbf{Z} \rangle], \quad (94b)$$

where the expectation is taken w.r.t.  $\{\mathbf{A}_i\}$  while assuming that  $(\mathbf{X}, \mathbf{Z}, \tau)$  are fixed. With these preparations in place, we set out to prove that: if  $m \geq C_0 nr \delta^{-2} \log m$ , then with probability at least  $1 - Ce^{-cn}$ ,

$$|E_m^\chi(\mathbf{X}, \mathbf{Z}, \tau) - E^\chi(\mathbf{X}, \mathbf{Z}, \tau)| \leq \delta\tau/2 \quad (95)$$

holds simultaneously for all  $(\mathbf{X}, \mathbf{Z}, \tau) \in \mathcal{T}_{\text{TRIP}}$ ; here  $C_0 > 0$  is some sufficiently large constant, and  $C, c > 0$  are some universal constants.

First, consider any fixed point  $(\mathbf{X}, \mathbf{Z}, \tau) \in \mathcal{T}_{\text{TRIP}}$ . Note that  $|\langle \mathbf{A}_i, \mathbf{X} \rangle \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau)| \leq \tau$  is bounded, and that the subgaussian norm of  $\langle \mathbf{A}_i, \mathbf{Z} \rangle$  obeys  $\|\langle \mathbf{A}_i, \mathbf{Z} \rangle\|_{\psi_2} = \|\mathcal{N}(0, 1)\|_{\psi_2} \lesssim 1$ . As a result,

$$\|\langle \mathbf{A}_i, \mathbf{X} \rangle \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau) \langle \mathbf{A}_i, \mathbf{Z} \rangle - E^\chi(\mathbf{X}, \mathbf{Z}, \tau)\|_{\psi_2} \lesssim \tau.$$

Invoking [Ver18, Theorem 2.6.2] tells us that for all  $t \geq 0$ ,

$$\mathbb{P}\left(|E_m^\chi(\mathbf{X}, \mathbf{Z}, \tau) - E^\chi(\mathbf{X}, \mathbf{Z}, \tau)| \geq t\tau\right) \leq 2 \exp(-c_1 mt^2)$$

holds for some constant  $c_1 > 0$ . Next, we construct an  $\epsilon$ -net to cover  $\mathcal{T}_{\text{TRIP}}$ . In view of [CP11, Lemma 3.1], the set  $\mathcal{R}_r^{\text{norm}}$  defined in (90) has an  $\epsilon$ -net (in terms of  $\|\cdot\|_{\text{F}}$  distance) of cardinality at most  $(9/\epsilon)^{3nr}$ . In addition, we can cover the interval  $[m^{-100}, 1.35]$  with precision  $\epsilon_\tau$  using no more than  $2/\epsilon_\tau$  equidistant points. Putting all this together, we can construct a set  $\mathcal{M}_{\text{TRIP}} \subseteq \mathcal{R}_r^{\text{norm}} \times \mathcal{R}_r^{\text{norm}} \times [0, 1.35]$  of cardinality at most  $(9/\epsilon)^{6nr}(2/\epsilon_\tau)$  such that: for any  $(\mathbf{X}, \mathbf{Z}, \tau) \in \mathcal{T}_{\text{TRIP}}$ , there exists some point  $(\mathbf{X}_0, \mathbf{Z}_0, \tau_0) \in \mathcal{M}_{\text{TRIP}}$  obeying

$$\|\mathbf{X} - \mathbf{X}_0\|_{\text{F}} \leq \epsilon, \quad \|\mathbf{Z} - \mathbf{Z}_0\|_{\text{F}} \leq \epsilon, \quad \text{and} \quad \tau - \epsilon_\tau \leq \tau_0 \leq \tau. \quad (96)$$

The union bound then implies that with probability at least  $1 - 2 \exp(-c_1 mt^2)(9/\epsilon)^{6nr}(2/\epsilon_\tau)$ , one has

$$|E_m^\chi - E^\chi(\mathbf{X}, \mathbf{Z}, \tau)| \leq t\tau, \quad \text{for all } (\mathbf{X}, \mathbf{Z}, \tau) \in \mathcal{M}_{\text{TRIP}}. \quad (97)$$

In what follows, we shall choose

$$t = \frac{1}{4}\delta \quad \text{and} \quad m \geq C_3 \frac{1}{\delta^2} \left( nr \log \frac{9}{\epsilon} + \log \frac{2}{\epsilon_\tau} \right) \quad (98)$$

so as to achieve a uniformly small error  $t\tau = \delta\tau/4$  in (97) with probability at least  $1 - 2 \exp(-c_3 m\delta^2)$  for some universal constant  $c_3 > 0$ .

Now, for any  $(\mathbf{X}, \mathbf{Z}, \tau) \in \mathcal{T}_{\text{TRIP}}$ , let  $(\mathbf{X}_0, \mathbf{Z}_0, \tau_0) \in \mathcal{M}_{\text{TRIP}}$  be the point satisfying (96). Then we have

$$|E_m^\chi(\mathbf{X}, \mathbf{Z}, \tau) - E^\chi(\mathbf{X}, \mathbf{Z}, \tau)| \leq \underbrace{|E_m^\chi(\mathbf{X}_0, \mathbf{Z}_0, \tau_0) - E^\chi(\mathbf{X}_0, \mathbf{Z}_0, \tau_0)|}_{(A)} \quad (99a)$$

$$+ \underbrace{|E_m^\chi(\mathbf{X}, \mathbf{Z}, \tau) - E_m^\chi(\mathbf{X}_0, \mathbf{Z}_0, \tau_0)|}_{(B)} + \underbrace{|E^\chi(\mathbf{X}, \mathbf{Z}, \tau) - E^\chi(\mathbf{X}_0, \mathbf{Z}_0, \tau_0)|}_{(C)}. \quad (99b)$$

Here, (A) is already bounded by  $\delta\tau/4$  by construction. We can control (B) via the following decomposition:

$$(B) \leq \underbrace{\left| \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X} \rangle \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau) \langle \mathbf{A}_i, \mathbf{Z} - \mathbf{Z}_0 \rangle \right|}_{(B.1)}$$

$$\begin{aligned}
& + \underbrace{\left| \frac{1}{m} \sum_{i=1}^m \left( \langle \mathbf{A}_i, \mathbf{X} \rangle \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau) - \langle \mathbf{A}_i, \mathbf{X}_0 \rangle \chi(\langle \mathbf{A}_i, \mathbf{X}_0 \rangle; \tau) \right) \langle \mathbf{A}_i, \mathbf{Z}_0 \rangle \right|}_{\text{(B.2)}} \\
& + \underbrace{\left| \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X}_0 \rangle \left( \chi(\langle \mathbf{A}_i, \mathbf{X}_0 \rangle; \tau) - \chi(\langle \mathbf{A}_i, \mathbf{X}_0 \rangle; \tau_0) \right) \langle \mathbf{A}_i, \mathbf{Z}_0 \rangle \right|}_{\text{(B.3)}}.
\end{aligned}$$

In light of the  $(2r, \delta)$ -RIP, the aforementioned properties of  $\chi$ , and the Cauchy-Schwarz inequality, we have

$$(B.1) \stackrel{(i)}{\leq} \tau \frac{1}{m} \sum_{i=1}^m |\langle \mathbf{A}_i, \mathbf{Z} - \mathbf{Z}_0 \rangle| \leq \tau \sqrt{\frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{Z} - \mathbf{Z}_0 \rangle^2} \lesssim \tau \epsilon,$$

$$(B.2) \leq \frac{1}{m} \sum_{i=1}^m |\langle \mathbf{A}_i, \mathbf{X} \rangle \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau) - \langle \mathbf{A}_i, \mathbf{X}_0 \rangle \chi(\langle \mathbf{A}_i, \mathbf{X}_0 \rangle; \tau)| \cdot |\langle \mathbf{A}_i, \mathbf{Z}_0 \rangle|$$

$$\stackrel{(ii)}{\leq} \frac{1}{c_\chi} \frac{1}{m} \sum_{i=1}^m |\langle \mathbf{A}_i, \mathbf{X} - \mathbf{X}_0 \rangle| \cdot |\langle \mathbf{A}_i, \mathbf{Z}_0 \rangle| \leq \frac{1}{c_\chi} \sqrt{\frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X} - \mathbf{X}_0 \rangle^2} \sqrt{\frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{Z}_0 \rangle^2} \lesssim \frac{\epsilon}{c_\chi},$$

$$(B.3) \leq \|\chi(\cdot; \tau) - \chi(\cdot; \tau_0)\|_\infty \frac{1}{m} \sum_{i=1}^m |\langle \mathbf{A}_i, \mathbf{X}_0 \rangle| \cdot |\langle \mathbf{A}_i, \mathbf{Z}_0 \rangle|$$

$$\stackrel{(iii)}{\lesssim} \frac{\epsilon_\tau}{c_\chi \tau} \sqrt{\frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X}_0 \rangle^2} \sqrt{\frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{Z}_0 \rangle^2} \lesssim \frac{\epsilon_\tau}{c_\chi \tau}.$$

Here, (i) uses  $|\langle \mathbf{A}_i, \mathbf{X} \rangle \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau)| \leq \tau$ , (ii) follows from the property that the function  $f(a) = a \cdot \chi(a; \tau)$  is  $1/c_\chi$ -Lipschitz continuous, whereas (iii) is due to the property  $\|\chi(\cdot; \tau) - \chi(\cdot; \tau_0)\|_\infty \leq \epsilon_\tau / (c_\chi \tau)$ . The term (C) can be controlled by the same decomposition and thus enjoys the same upper bound. Putting these back into (99), we have for all  $(\mathbf{X}, \mathbf{Z}, \tau) \in \mathcal{T}_{\text{TRIP}}$ ,

$$|E_m^\chi(\mathbf{X}, \mathbf{Z}, \tau) - E^\chi(\mathbf{X}, \mathbf{Z}, \tau)| \leq \frac{1}{4} \delta \tau + C_3 \left( \tau \epsilon + \frac{\epsilon}{c_\chi} + \frac{\epsilon_\tau}{c_\chi \tau} \right)$$

for some universal constant  $C_3 > 0$ . Recalling that  $\tau \geq m^{-100}$ , and choosing  $\epsilon \leq c_4 \delta c_\chi m^{-100}$  and  $\epsilon_\tau \leq c_5 \delta c_\chi m^{-200}$  for some sufficiently small constants  $c_4, c_5 > 0$ , we have

$$|E_m^\chi(\mathbf{X}, \mathbf{Z}, \tau) - E^\chi(\mathbf{X}, \mathbf{Z}, \tau)| \leq \delta \tau / 2.$$

Plugging our choice of  $\epsilon$  and  $\epsilon_\tau$  into (98) immediately establishes the claim (95) of this step.

**Step 2: controlling the errors incurred by using the surrogate  $\chi$ .** Similar to (94), we define

$$\begin{aligned}
E_m(\mathbf{X}, \mathbf{Z}, \tau) & := \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X} \rangle \mathbb{1}(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau) \langle \mathbf{A}_i, \mathbf{Z} \rangle, \\
E(\mathbf{X}, \mathbf{Z}, \tau) & := \mathbb{E}[\langle \mathbf{A}_i, \mathbf{X} \rangle \mathbb{1}(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau) \langle \mathbf{A}_i, \mathbf{Z} \rangle] = w(\tau) \langle \mathbf{X}, \mathbf{Z} \rangle,
\end{aligned}$$

where the expectation is taken assuming independence between  $\mathbf{A}_i$  and  $(\mathbf{X}, \mathbf{Z}, \tau)$ . In this step, we aim to show that: if  $m \geq C_0 n r \delta^{-2} \log m$ , then with probability at least  $1 - Ce^{-cn}$ ,

$$|E_m(\mathbf{X}, \mathbf{Z}, \tau) - E(\mathbf{X}, \mathbf{Z}, \tau)| \leq |E_m^\chi(\mathbf{X}, \mathbf{Z}, \tau) - E^\chi(\mathbf{X}, \mathbf{Z}, \tau)| + \delta \tau / 2 \quad (100)$$

holds simultaneously for all  $(\mathbf{X}, \mathbf{Z}, \tau) \in \mathcal{T}_{\text{TRIP}}$ . If this were true, then combining this with (95) would immediately conclude the proof of Lemma 1.

Towards establishing (100), we start with the following decomposition:

$$\begin{aligned}
|E_m(\mathbf{X}, \mathbf{Z}, \tau) - E(\mathbf{X}, \mathbf{Z}, \tau)| &\leq |E_m^\chi(\mathbf{X}, \mathbf{Z}, \tau) - E^\chi(\mathbf{X}, \mathbf{Z}, \tau)| + \underbrace{|E(\mathbf{X}, \mathbf{Z}, \tau) - E^\chi(\mathbf{X}, \mathbf{Z}, \tau)|}_{(A)} \\
&\quad + \underbrace{|E_m(\mathbf{X}, \mathbf{Z}, \tau) - E_m^\chi(\mathbf{X}, \mathbf{Z}, \tau)|}_{(B)}, \tag{101}
\end{aligned}$$

where we abuse the notation (A) and (B). In the sequel, we shall control (A) and (B) separately.

- Regarding (A), the Cauchy-Schwarz inequality gives

$$\begin{aligned}
(A) &= \left| \mathbb{E}[\langle \mathbf{A}_i, \mathbf{X} \rangle (\mathbb{1}(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau) - \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau)) \langle \mathbf{A}_i, \mathbf{Z} \rangle] \right| \\
&\leq \sqrt{\mathbb{E}[(\mathbb{1}(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau) - \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau))^2]} \sqrt{\mathbb{E}[(\langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle)^2]} \lesssim \sqrt{c_\chi \tau}.
\end{aligned}$$

The last inequality holds since  $|\mathbb{1}(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau) - \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau)| \in [0, 1]$  is non-zero only for  $|\langle \mathbf{A}_i, \mathbf{X} \rangle|$  on an interval of length  $c_\chi \tau$ , over which the probability density function of  $\langle \mathbf{A}_i, \mathbf{X} \rangle \sim \mathcal{N}(0, 1)$  is upper bounded by some constant. By our choice of  $c_\chi$  in (93), we have (A)  $\leq \delta \tau / 4$ .

- We then move on to (B). For notational convenience, given any  $\tau > 0$ , we let

$$\tau' = \tau'(\tau) := \frac{\tau}{1 - c_\chi}, \tag{102}$$

which clearly satisfies  $\chi(a; \tau) \leq \mathbb{1}(a; \tau) \leq \chi(a; \tau')$ . In addition, defining

$$\mathbb{1}_-(a) := \mathbb{1}(a < 0), \quad \mathbb{1}_+(a) := \mathbb{1}(a \geq 0), \quad a \in \mathbb{R},$$

we can deduce that

$$\begin{aligned}
E_m(\mathbf{X}, \mathbf{Z}, \tau) &\leq E_m^\chi(\mathbf{X}, \mathbf{Z}, \tau) + \frac{1}{m} \sum_{i=1}^m (\chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau') - \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau)) \langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle \mathbb{1}_+(\langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle), \\
E_m(\mathbf{X}, \mathbf{Z}, \tau) &\geq E_m^\chi(\mathbf{X}, \mathbf{Z}, \tau) + \frac{1}{m} \sum_{i=1}^m (\chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau') - \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau)) \langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle \mathbb{1}_-(\langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle).
\end{aligned}$$

As a consequence,

$$\begin{aligned}
(B) &\leq \max \left\{ \underbrace{\left| \frac{1}{m} \sum_{i=1}^m (\chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau') - \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau)) \langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle \mathbb{1}_+(\langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle) \right|}_{(C)}, \right. \\
&\quad \left. \left| \frac{1}{m} \sum_{i=1}^m (\chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau') - \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau)) \langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle \mathbb{1}_-(\langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle) \right| \right\}.
\end{aligned}$$

Next, we demonstrate how to analyze the first term (C) above; the analysis for the other term is essentially the same. For notational simplicity, define

$$F_m^+(\mathbf{X}, \mathbf{Z}, \tau) := \frac{1}{m} \sum_{i=1}^m \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau) \langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle \mathbb{1}_+(\langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle),$$

$$E^+(\mathbf{X}, \mathbf{Z}, \tau) := \mathbb{E} \left[ (\chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau') - \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau)) \langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle \mathbb{1}_+(\langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle) \right],$$

where the expectation is again taken assuming that  $\mathbf{A}_i$  is independent of  $\mathbf{X}, \mathbf{Z}$  and  $\tau$ . Then we have

$$(C) = |F_m^+(\mathbf{X}, \mathbf{Z}, \tau') - F_m^+(\mathbf{X}, \mathbf{Z}, \tau)| \leq |E^+(\mathbf{X}, \mathbf{Z}, \tau)| + |F_m^+(\mathbf{X}, \mathbf{Z}, \tau') - F_m^+(\mathbf{X}, \mathbf{Z}, \tau) - E^+(\mathbf{X}, \mathbf{Z}, \tau)|.$$

Regarding the first term on the right-hand side, we follow an argument similar to our previous analysis for (A) to obtain

$$|E^+(\mathbf{X}, \mathbf{Z}, \tau)| \leq \sqrt{\mathbb{E}\left[\left(\chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau') - \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau)\right)^2\right]} \sqrt{\mathbb{E}\left[\left(\langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle\right)^2\right]} \lesssim \sqrt{c_\chi \tau} \leq c_2 \delta \tau$$

for some sufficiently small constant  $0 < c_2 < 1/8$ . Thus, it remains to show that

$$\left|F_m^+(\mathbf{X}, \mathbf{Z}, \tau') - F_m^+(\mathbf{X}, \mathbf{Z}, \tau) - E^+(\mathbf{X}, \mathbf{Z}, \tau)\right| \leq \frac{1}{8} \delta \tau \quad (103)$$

holds simultaneously for all  $(\mathbf{X}, \mathbf{Z}, \tau) \in \mathcal{T}_{\text{TRIP}}$ . Note that by definition,  $F_m^+(\mathbf{X}, \mathbf{Z}, \tau)$  is the empirical average of some Lipschitz continuous function (in particular,  $\langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle \mathbb{1}_+(\langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle)$  is 1-Lipschitz continuous over  $\langle \mathbf{A}_i, \mathbf{X} \rangle \langle \mathbf{A}_i, \mathbf{Z} \rangle$ ). Therefore, we can prove (103) by a standard covering argument similar to that in Step 1; we omit the details for brevity. Putting the above bounds together, we establish that (B)  $\leq \delta \tau / 4$ .

- Combining the above bounds (A)  $\leq \delta \tau / 4$  and (B)  $\leq \delta \tau / 4$  with (101), we finish the proof of (100).

**Proof for the case  $0 \leq \tau < m^{-100}$ .** It remains to prove that (91) holds simultaneously for all  $\mathbf{X}, \mathbf{Z} \in \mathcal{R}_r^{\text{norm}}$  (cf. (90)) and all  $0 \leq \tau < m^{-100}$ . We start with the following decomposition:

$$\left| \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X} \rangle \mathbb{1}(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau) \langle \mathbf{A}_i, \mathbf{Z} \rangle - w(\tau) \langle \mathbf{X}, \mathbf{Z} \rangle \right| \leq \left| \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X} \rangle \mathbb{1}(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau) \langle \mathbf{A}_i, \mathbf{Z} \rangle \right| + |w(\tau) \langle \mathbf{X}, \mathbf{Z} \rangle|. \quad (104)$$

The second term on the right-hand side of (104) can be bounded by

$$|w(\tau) \langle \mathbf{X}, \mathbf{Z} \rangle| \leq w(\tau) \stackrel{(i)}{\leq} \tau^3 \leq m^{-200} \tau \stackrel{(ii)}{\leq} 0.1 \delta \tau,$$

where (i) can be seen from the definition of  $w(\cdot)$  in (60), and (ii) relies on the observation that our assumption  $m \geq C_0 n r \delta^{-2} \log m$  implies  $\delta \gtrsim m^{-1/2}$ .

It thus remains to show that the first term on the right-hand side of (104) is bounded by  $0.9 \delta \tau$ . In view of  $(2r, \delta)$ -RIP, the Cauchy-Schwarz inequality, and the observation that  $|\langle \mathbf{A}_i, \mathbf{X} \rangle \mathbb{1}(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau)| \leq \tau$ , we have

$$\begin{aligned} \left| \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X} \rangle \mathbb{1}(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau) \langle \mathbf{A}_i, \mathbf{Z} \rangle \right| &\leq \sqrt{\frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X} \rangle^2 \mathbb{1}(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau)} \cdot \sqrt{\frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{Z} \rangle^2} \\ &\leq 2 \sqrt{\frac{1}{m} \sum_{i=1}^m \tau^2 \mathbb{1}(\langle \mathbf{A}_i, \mathbf{X} \rangle; \tau)} \leq 2\tau \sqrt{\frac{1}{m} \sum_{i=1}^m \mathbb{1}(\langle \mathbf{A}_i, \mathbf{X} \rangle; m^{-100})}, \end{aligned} \quad (105)$$

where the last inequality uses the assumption that  $\tau < m^{-100}$ . We can invoke Lemma 3 with  $t = 0.01 \delta^2$  and  $\epsilon = m^{-200}$  to obtain that with probability at least  $1 - Ce^{-cn}$  (for some constants  $c, C > 0$ ),

$$\frac{1}{m} \sum_{i=1}^m \mathbb{1}(\langle \mathbf{A}_i, \mathbf{X} \rangle; m^{-100}) \leq \mathbb{P}(Z_0 \leq 1.01 m^{-100}) + t + \frac{200\epsilon}{m^{-100}} \leq 2t = 0.02 \delta^2$$

holds simultaneously for all  $\mathbf{X} \in \mathcal{R}_r^{\text{norm}}$ , provided that  $m \geq C_0 n r \delta^{-2} \log m$ ; here,  $Z_0$  denotes a random variable having the same distribution as  $|\mathcal{N}(0, 1)|$ . Plugging this into (105) confirms that the first term on the right-hand side of (104) is bounded by  $0.9 \delta \tau$ , thus concluding the proof for the case with  $0 \leq \tau < m^{-100}$ .

## C.2 Proof of Lemma 2

It is easy to check that  $Q_\alpha(\mathcal{D}) / \|\mathbf{X}_1\|_{\text{F}}$  is invariant under a global scaling of  $\{\mathbf{X}_1, \dots, \mathbf{X}_K\}$ . Therefore, it suffices to consider a normalized version of  $\mathcal{T}_1$  (86) defined as follows:

$$\mathcal{T}_1^{\text{norm}} := \mathcal{T}_1 \cap \{(\mathbf{X}_1, \dots, \mathbf{X}_k) : \|\mathbf{X}_1\|_{\text{F}} = 1\}. \quad (106)$$

In what follows, we shall treat the upper bound and the lower bound separately and invoke a standard covering argument to prove Lemma 2 with  $\mathcal{T}_1$  replaced by  $\mathcal{T}_1^{\text{norm}}$ . Throughout this proof, we denote by  $Z$  a random variable following the distribution of  $|\mathcal{N}(0, 1)|$ .

**Step 1: upper bounding  $Q_\alpha(\mathcal{D})$ .** Since  $\alpha \leq 0.8p_1$ , we have

$$Q_\alpha(\mathcal{D}) \leq Q_{\alpha/p_1}(\mathcal{D}_1) \leq Q_{0.8}(\mathcal{D}_1).$$

Now it suffices to upper bound  $Q_{0.8}(\mathcal{D}_1)$ , which is only related to  $\mathbf{X}_1 \in \mathcal{R}_r^{\text{norm}}$ . Consider any fixed point  $\mathbf{X}_1 \in \mathcal{R}_r^{\text{norm}}$ . Note that the set  $\mathcal{D}_1$  defined in (85) contains i.i.d. samples having the same distribution as  $Z$ . This combined with the concentration of empirical quantiles [Ser09, Section 2.3.2] gives

$$\mathbb{P}(Q_{0.8}(\mathcal{D}_1) \geq Q_{0.8}(Z) + 0.01) \leq \exp(-c_2 N_1) \quad (107)$$

for some universal constant  $c_2 > 0$ . Here,  $N_1 := |\Omega_1^*| \asymp N/K$  by the assumption of the well-balancedness property (18). Next, we construct an  $\epsilon$ -net of  $\mathcal{R}_r^{\text{norm}}$  — denoted by  $\mathcal{M}$  — whose cardinality is at most  $(9/\epsilon)^{3nr}$  (according to [CP11, Lemma 3.1]). Taking the union bound over  $\mathcal{M}$  and assuming that

$$N_1 \geq C_0 nr \log \frac{9}{\epsilon}$$

for some sufficiently large constant  $C_0 > 0$ , we have with probability at least  $1 - Ce^{-cn}$ , for all  $\mathbf{X}_1 \in \mathcal{M}$ , the dataset  $\mathcal{D}_1$  defined in (85) satisfies  $Q_{0.8}(\mathcal{D}_1) \leq Q_{0.8}(Z) + 0.01$ . Finally, consider an arbitrary  $\mathbf{X}_1 \in \mathcal{R}_r^{\text{norm}}$ , and let  $\mathbf{X}_1^0$  be the point in  $\mathcal{M}$  such that  $\|\mathbf{X}_1^0 - \mathbf{X}_1\|_F \leq \epsilon$ . Denote by  $\mathcal{D}_1^0$  the dataset generated by  $\mathbf{X}_1^0$  analogous to (85). Then we have

$$|Q_{0.8}(\mathcal{D}_1) - Q_{0.8}(\mathcal{D}_1^0)| \leq \max_{i \in \Omega_1^*} |\langle \mathbf{A}_i, \mathbf{X}_1 \rangle - \langle \mathbf{A}_i, \mathbf{X}_1^0 \rangle| \leq \epsilon \max_{i \in \Omega_1^*} \|\mathbf{A}_i\|_F \lesssim \epsilon (n + \sqrt{\log N_1}),$$

where the last inequality holds with probability at least  $1 - Ce^{-cn}$ , according to Proposition 2. Setting  $\epsilon = N_1^{-10}$ , we further have  $|Q_{0.8}(\mathcal{D}_1) - Q_{0.8}(\mathcal{D}_1^0)| \lesssim N_1^{-9} \leq 0.01$ , as long as  $N_1$  is sufficiently large. In addition, it can be verified numerically that  $Q_{0.8}(Z) < 1.30$ . These together imply that for any  $(\mathbf{X}_1, \dots, \mathbf{X}_K) \in \mathcal{T}_1^{\text{norm}}$ , we have

$$Q_\alpha(\mathcal{D}) \leq Q_{0.8}(\mathcal{D}_1) \leq Q_{0.8}(Z) + 0.02 \leq 1.35,$$

which gives rise to the upper bound in Lemma 2.

**Step 2: lower bounding  $Q_\alpha(\mathcal{D})$ .** For notational convenience, we denote

$$q := \frac{0.7\alpha}{p_1} \in [0.42, 0.56], \quad \text{and} \quad B_N := \frac{1}{N} \sum_{k=1}^K \sum_{i \in \Omega_k^*} \mathbb{1}\left(|\langle \mathbf{A}_i, \mathbf{X}_k \rangle| \leq \frac{Q_q(Z)}{1.01}\right). \quad (108)$$

Clearly, by the definition of  $B_N$ , one has

$$\mathbb{P}\left(Q_\alpha(\mathcal{D}) < \frac{Q_q(Z)}{1.01}\right) \leq \mathbb{P}(B_N > \alpha),$$

where it can be verified numerically that  $Q_q(Z)/1.01 \geq 0.54$ . Therefore, it suffices to upper bound the probability  $\mathbb{P}(B_N > \alpha)$ . To accomplish this, we first upper bound  $B_N$  as follows:

$$\begin{aligned} B_N &= \frac{1}{N} \sum_{k=1}^K \sum_{i \in \Omega_k^*} \mathbb{1}\left(\left|\left\langle \mathbf{A}_i, \frac{\mathbf{X}_k}{\|\mathbf{X}_k\|_F} \right\rangle\right| \leq \frac{Q_q(Z)}{1.01\|\mathbf{X}_k\|_F}\right) \\ &\leq \frac{1}{N} \sum_{i \in \Omega_1^*} \mathbb{1}\left(|\langle \mathbf{A}_i, \mathbf{X}_1 \rangle| \leq \frac{Q_q(Z)}{1.01}\right) + \frac{1}{N} \sum_{k \neq 1} \sum_{i \in \Omega_k^*} \mathbb{1}\left(\left|\left\langle \mathbf{A}_i, \frac{\mathbf{X}_k}{\|\mathbf{X}_k\|_F} \right\rangle\right| \leq \frac{c_0 Q_q(Z)}{1.01K}\right). \end{aligned} \quad (109)$$

Here, the last line follows from the assumption that  $1 = \|\mathbf{X}_1\|_F \leq (c_0/K) \min_{k \neq 1} \|\mathbf{X}_k\|_F$ ; see the definition of  $\mathcal{T}_1^{\text{norm}}$  in (106). Note that  $\mathbf{X}_1 \in \mathcal{R}_r^{\text{norm}}$ , and for all  $k \neq 1$ , we also have  $\mathbf{X}_k/\|\mathbf{X}_k\|_F \in \mathcal{R}_r^{\text{norm}}$ . Therefore, we can invoke Lemma 3 with  $m = N_1 = |\Omega_1^*|$ ,  $\tau = Q_q(Z)/1.01$ ,  $t = 0.15\alpha$  and  $\epsilon = N_1^{-10}$  to obtain that: with probability at least  $1 - Ce^{-cn}$  (provided that  $m \geq C_0nrK^2 \log m$ ), the following holds simultaneously for all  $\mathbf{X}_1 \in \mathcal{R}_r^{\text{norm}}$ :

$$\frac{1}{N_1} \sum_{i \in \Omega_1^*} \mathbb{1}\left(\left|\langle \mathbf{A}_i, \mathbf{X}_1 \rangle\right| \leq \frac{Q_q(Z)}{1.01}\right) \leq \mathbb{P}(Z \leq Q_q(Z)) + t + \frac{200\epsilon}{\tau} = q + 0.15\alpha + \frac{202N_1^{-10}}{Q_q(Z)}.$$

Similarly, for all  $k \neq 1$ , one can apply Lemma 3 with  $m = N_k := |\Omega_k^*|$ ,  $\tau = c_0Q_q(Z)/(1.01K)$ ,  $t = 0.15\alpha$  and  $\epsilon = N_k^{-10}$  to show that: with probability at least  $1 - Ce^{-cn}$  (provided  $m \geq C_0nrK^2 \log m$ ), the following holds simultaneously for all  $\mathbf{X}_k/\|\mathbf{X}_k\|_F \in \mathcal{R}_r^{\text{norm}}$ :

$$\frac{1}{N_k} \sum_{i \in \Omega_k^*} \mathbb{1}\left(\left|\left\langle \mathbf{A}_i, \frac{\mathbf{X}_k}{\|\mathbf{X}_k\|_F} \right\rangle\right| \leq \frac{c_0Q_q(Z)}{1.01K}\right) \leq \mathbb{P}\left(Z \leq \frac{c_0Q_q(Z)}{K}\right) + t + \frac{200\epsilon}{\tau} \leq \frac{c_0Q_q(Z)}{K} + 0.15\alpha + \frac{202KN_k^{-10}}{c_0Q_q(Z)},$$

where the last inequality relies on the property of  $Z$ . Combine the above two bounds with (109) to reach

$$\begin{aligned} B_N &\leq p_1 \left( q + 0.15\alpha + \frac{202N_1^{-10}}{Q_q(Z)} \right) + \sum_{k \neq 1} p_k \left( \frac{c_0Q_q(Z)}{K} + 0.15\alpha + \frac{202KN_k^{-10}}{c_0Q_q(Z)} \right) \\ &\leq p_1q + \frac{c_0Q_q(Z)}{K} + 0.15\alpha + p_1 \frac{202N_1^{-10}}{Q_q(Z)} + \sum_{k \neq 1} p_k \frac{202KN_k^{-10}}{c_0Q_q(Z)}. \end{aligned}$$

Recall that  $p_1q = 0.7\alpha$ ,  $\alpha \asymp p_1 \asymp 1/K$ , and observe that  $p_1 \frac{202N_1^{-10}}{Q_q(Z)} + \sum_{k \neq 1} p_k \frac{202KN_k^{-10}}{c_0Q_q(Z)} \leq 0.05\alpha$  as long as  $N_k \gtrsim K$  for all  $k$ . Putting these together guarantees that  $B_N \leq \alpha$  as desired, which further implies

$$Q_\alpha(\mathcal{D}) \geq Q_q(Z)/1.01 \geq 0.54.$$

Combining this lower bound with the upper bound in Step 1 completes our proof of Lemma 2.

### C.3 Proof of Lemma 3

Throughout the proof, we assume that the ensemble  $\{\mathbf{A}_i\}$  obeys  $(2r, 1/4)$ -RIP. In view of Lemma 1, this happens with probability at least  $1 - C_2e^{-c_2n}$  for some constants  $c_2, C_2 > 0$ , as long as  $m \geq Cnr$ . Recall the definition of  $\chi$  from Appendix C.1, and set the parameter as  $c_\chi = 0.01/1.01$ . One then has

$$\frac{1}{m} \sum_{i=1}^m \mathbb{1}(|\langle \mathbf{A}_i, \mathbf{X} \rangle| \leq \tau) \leq \frac{1}{m} \sum_{i=1}^m \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; 1.01\tau)$$

In the sequel, we invoke the standard covering argument to upper bound  $\frac{1}{m} \sum_{i=1}^m \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; 1.01\tau)$ .

First, consider a fixed  $\mathbf{X} \in \mathcal{R}_r^{\text{norm}}$  independent of  $\{\mathbf{A}_i\}$ . In this case we can bound the expectation as

$$\mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; 1.01\tau) \right] \leq \mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \mathbb{1}(|\langle \mathbf{A}_i, \mathbf{X} \rangle| \leq 1.01\tau) \right] = \mathbb{P}(Z \leq 1.01\tau),$$

where we recall that  $Z$  follows the same distribution as  $|\mathcal{N}(0, 1)|$ . In addition, note that  $\frac{1}{m} \sum_{i=1}^m \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; 1.01\tau)$  is the empirical average of  $m$  independent random variables, each lying within  $[0, 1]$  and having variance bounded by  $2\tau$ . Therefore, for all  $t \geq 0$ , one sees from Bernstein's inequality [Ver18, Theorem 2.8.4] that

$$\mathbb{P} \left( \frac{1}{m} \sum_{i=1}^m \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; 1.01\tau) \geq \mathbb{E} \left[ \frac{1}{m} \sum_{i=1}^m \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; 1.01\tau) \right] + t \right) \leq \exp \left( -\frac{c_1mt^2}{\tau + t} \right),$$



where  $c_0, c_1 > 0$  are some universal constants. Let  $\mathcal{M} \subseteq \mathcal{R}_r^{\text{norm}}$  be an  $\epsilon$ -net of  $\mathcal{R}_r^{\text{norm}}$ , whose cardinality is at most  $(9/\epsilon)^{3nr}$ . The union bound reveals that: with probability at least  $1 - (9/\epsilon)^{3nr} \exp(-c_1 m t^2 / (\tau + t))$ , one has

$$\sup_{\mathbf{X} \in \mathcal{M}} \frac{1}{m} \sum_{i=1}^m \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; 1.01\tau) \leq \mathbb{P}(Z \leq 1.01\tau) + t.$$

Next, we move on to account for an arbitrary  $\mathbf{X} \in \mathcal{R}_r^{\text{norm}}$  (which is not necessarily independent of  $\{\mathbf{A}_i\}$ ). Let  $\mathbf{X}_0$  be a point in  $\mathcal{M}$  obeying  $\|\mathbf{X} - \mathbf{X}_0\|_{\text{F}} \leq \epsilon$ . As a result, one has

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m \chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; 1.01\tau) - \frac{1}{m} \sum_{i=1}^m \chi(\langle \mathbf{A}_i, \mathbf{X}_0 \rangle; 1.01\tau) &\leq \frac{1}{m} \sum_{i=1}^m |\chi(\langle \mathbf{A}_i, \mathbf{X} \rangle; 1.01\tau) - \chi(\langle \mathbf{A}_i, \mathbf{X}_0 \rangle; 1.01\tau)| \\ &\stackrel{\text{(i)}}{\leq} \frac{100}{\tau} \cdot \frac{1}{m} \sum_{i=1}^m |\langle \mathbf{A}_i, \mathbf{X} - \mathbf{X}_0 \rangle| \\ &\stackrel{\text{(ii)}}{\leq} \frac{100}{\tau} \cdot \sqrt{\frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X} - \mathbf{X}_0 \rangle^2} \\ &\stackrel{\text{(iii)}}{\leq} \frac{200}{\tau} \|\mathbf{X} - \mathbf{X}_0\|_{\text{F}} \leq \frac{200}{\tau} \epsilon. \end{aligned}$$

Here the inequality (i) holds since  $\chi(\cdot; 1.01\tau)$  is Lipschitz with the Lipschitz constant  $1/(1.01c_\chi\tau) = 100/\tau$ , the relation (ii) results from the Cauchy-Schwarz inequality, and (iii) follows since  $\{\mathbf{A}_i\}$  obeys  $(2r, 1/4)$ -RIP.

Combine the above two inequalities to finish the proof.

## D Estimating unknown parameters in Algorithm 4

Throughout the paper, we have assumed the knowledge of several problem-specific parameters, e.g. the proportion  $p_k$  of the  $k$ -th component, the rank  $r_k$  of the low-rank matrix  $\mathbf{M}_k^*$  and the rank  $R = \text{rank}(\mathbb{E}[\mathbf{Y}])$ . In the sequel, we specify where we need them and discuss how to estimate them in practice.

- In Line 2 of Algorithm 4, when running Algorithm 1, we need to know  $R = \text{rank}(\mathbb{E}[\mathbf{Y}])$ , which can be estimated faithfully by examining the singular values of the data matrix  $\mathbf{Y}$ .
- In Line 3 of Algorithm 4, when running Algorithm 2, we need to know  $\{r_k\}_{1 \leq k \leq K}$ , where  $r_k = \text{rank}(\mathbf{M}_k^*)$ . Recall from (14) that  $\mathbf{U}\widehat{\mathbf{S}}_k\mathbf{V}^\top \approx \mathbf{M}_k^*$ ; therefore,  $r_k$  can be estimated accurately by examining the singular values of  $\widehat{\mathbf{S}}_k$ .
- In Line 5 of Algorithm 4, when running Algorithm 3, we need to know  $p_k$  to set  $\eta_k$  and  $\alpha_k$  appropriately. It turns out that the outputs  $\{\omega_k\}$  of the tensor method (see Algorithm 5) satisfy  $\omega_k \approx p_k, 1 \leq k \leq K$ .

## References

- [ACHL19] S. Arora, N. Cohen, W. Hu, and Y. Luo. Implicit regularization in deep matrix factorization. In *Advances in Neural Information Processing Systems*, pages 7413–7424, 2019.
- [AEP07] A. Argyriou, T. Evgeniou, and M. Pontil. Multi-task feature learning. In *Advances in Neural Information Processing Systems*, pages 41–48, 2007.
- [AGH<sup>+</sup>14] A. Anandkumar, R. Ge, D. Hsu, S. M. Kakade, and M. Telgarsky. Tensor decompositions for learning latent variable models. *Journal of Machine Learning Research*, 15:2773–2832, 2014.
- [AZ05] R. K. Ando and T. Zhang. A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6(Nov):1817–1853, 2005.
- [Bax00] J. Baxter. A model of inductive bias learning. *Journal of Artificial Intelligence Research*, 12:149–198, 2000.

- [BDS03] S. Ben-David and R. Schuller. Exploiting task relatedness for multiple task learning. In *Learning Theory and Kernel Machines*, pages 567–580. Springer, 2003.
- [BJK15] K. Bhatia, P. Jain, and P. Kar. Robust regression via hard thresholding. In *Advances in Neural Information Processing Systems*, pages 721–729, 2015.
- [BNS16] S. Bhojanapalli, B. Neyshabur, and N. Srebro. Global optimality of local search for low rank matrix recovery. In *Advances in Neural Information Processing Systems*, pages 3873–3881, 2016.
- [Car97] R. Caruana. Multitask learning. *Machine Learning*, 28(1):41–75, 1997.
- [CC17] Y. Chen and E. J. Candès. Solving random quadratic systems of equations is nearly as easy as solving linear systems. *Communications on Pure and Applied Mathematics*, 70(5):822–883, 2017.
- [CC18a] Y. Chen and E. J. Candès. The projected power method: An efficient algorithm for joint alignment from pairwise differences. *Communications on Pure and Applied Mathematics*, 71(8):1648–1714, 2018.
- [CC18b] Y. Chen and Y. Chi. Harnessing structures in big data via guaranteed low-rank matrix estimation: Recent theory and fast algorithms via convex and nonconvex optimization. *IEEE Signal Processing Magazine*, 35(4):14–31, 2018.
- [CCD<sup>+</sup>19] V. Charisopoulos, Y. Chen, D. Davis, M. Díaz, L. Ding, and D. Drusvyatskiy. Low-rank matrix recovery with composite optimization: good conditioning and rapid convergence. *arXiv preprint arXiv:1904.10020*, 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, 2019.
- [CCF<sup>+</sup>ar] Y. Chen, Y. Chi, J. Fan, C. Ma, and Y. Yan. Noisy matrix completion: Understanding statistical guarantees for convex relaxation via nonconvex optimization. *SIAM Journal on Optimization*, to appear.
- [CCG15] Y. Chen, Y. Chi, and A. J. Goldsmith. Exact and stable covariance estimation from quadratic sampling via convex programming. *IEEE Transactions on Information Theory*, 61(7):4034–4059, 2015.
- [CFMY19] Y. Chen, J. Fan, C. Ma, and Y. Yan. Inference and uncertainty quantification for noisy matrix completion. *Proceedings of the National Academy of Sciences*, 116(46):22931–22937, 2019.
- [CGH14] Y. Chen, L. Guibas, and Q. Huang. Near-optimal joint object matching via convex relaxation. In *Proceedings of the International Conference on Machine Learning*, pages 1269–1277, 2014.
- [CL13] A. T. Chaganty and P. Liang. Spectral experts for estimating mixtures of linear regressions. In *Proceedings of the International Conference on Machine Learning*, pages 1040–1048, 2013.
- [CLC19] Y. Chi, Y. M. Lu, and Y. Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing*, 67(20):5239–5269, 2019.
- [CLL20] J. Chen, D. Liu, and X. Li. Nonconvex rectangular matrix completion via gradient descent without  $\ell_{2,\infty}$  regularization. *IEEE Transactions on Information Theory*, 2020.
- [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.
- [CLS20] S. Chen, J. Li, and Z. Song. Learning mixtures of linear regressions in subexponential time via fourier moments. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, pages 587–600, 2020.

- [CP11] E. J. Candes and Y. Plan. Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements. *IEEE Transactions on Information Theory*, 57(4):2342–2359, 2011.
- [CR09] E. J. Candès and B. Recht. Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9(6):717, 2009.
- [CYC14] Y. Chen, X. Yi, and C. Caramanis. A convex formulation for mixed regression with two components: Minimax optimal rates. In *Proceedings of the Conference on Learning Theory*, pages 560–604, 2014.
- [DC20] L. Ding and Y. Chen. Leave-one-out approach for matrix completion: Primal and dual analysis. *IEEE Transactions on Information Theory*, 2020.
- [DH00] P. Deb and A. M. Holmes. Estimates of use and costs of behavioural health care: a comparison of standard and finite mixture models. *Health Economics*, 9(6):475–489, 2000.
- [DHK<sup>+</sup>20] S. S. Du, W. Hu, S. M. Kakade, J. D. Lee, and Q. Lei. Few-shot learning via learning the representation, provably. *arXiv preprint arXiv:2002.09434*, 2020.
- [DV89] R. D. De Veaux. Mixtures of linear regressions. *Computational Statistics & Data Analysis*, 8(3):227–245, 1989.
- [EMP05] T. Evgeniou, C. A. Micchelli, and M. Pontil. Learning multiple tasks with kernel methods. *Journal of Machine Learning Research*, 6(Apr):615–637, 2005.
- [FAL17] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, pages 1126–1135. JMLR.org, 2017.
- [GJZ17] R. Ge, C. Jin, and Y. Zheng. No spurious local minima in nonconvex low rank problems: A unified geometric analysis. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 1233–1242. PMLR, 2017.
- [GS99] S. Gaffney and P. Smyth. Trajectory clustering with mixtures of regression models. In *Proceedings of the Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 63–72, 1999.
- [HJ18] P. Hand and B. Joshi. A convex program for mixed linear regression with a recovery guarantee for well-separated data. *Information and Inference: A Journal of the IMA*, 7(3):563–579, 2018.
- [JNS13] P. Jain, P. Netrapalli, and S. Sanghavi. Low-rank matrix completion using alternating minimization. In *Proceedings of the Forty-fifth Annual ACM Symposium on Theory of Computing*, pages 665–674, 2013.
- [JSRR10] A. Jalali, S. Sanghavi, C. Ruan, and P. K. Ravikumar. A dirty model for multi-task learning. In *Advances in neural information processing systems*, pages 964–972, 2010.
- [KC07] A. Khalili and J. Chen. Variable selection in finite mixture of regression models. *Journal of the American Statistical Association*, 102(479):1025–1038, 2007.
- [KHC20] J. Kwon, N. Ho, and C. Caramanis. On the minimax optimality of the em algorithm for learning two-component mixed linear regression. *arXiv preprint arXiv:2006.02601*, 2020.
- [KMMP19] A. Krishnamurthy, A. Mazumdar, A. McGregor, and S. Pal. Sample complexity of learning mixture of sparse linear regressions. In *Advances in Neural Information Processing Systems*, pages 10532–10541, 2019.

- [KMO10] R. H. Keshavan, A. Montanari, and S. Oh. Matrix completion from a few entries. *IEEE Transactions on Information Theory*, 56(6):2980–2998, 2010.
- [KQC<sup>+</sup>19] J. Kwon, W. Qian, C. Caramanis, Y. Chen, and D. Davis. Global convergence of the em algorithm for mixtures of two component linear regression. In *Proceedings of the Conference on Learning Theory*, pages 2055–2110, 2019.
- [KSKO20] W. Kong, R. Somani, S. Kakade, and S. Oh. Robust meta-learning for mixed linear regression with small batches. *arXiv preprint arXiv:2006.09702*, 2020.
- [KSS<sup>+</sup>20] W. Kong, R. Somani, Z. Song, S. Kakade, and S. Oh. Meta-learning for mixed linear regression. *arXiv preprint arXiv:2002.08936*, 2020.
- [KYB19] J. M. Klusowski, D. Yang, and W. Brinda. Estimating the coefficients of a mixture of two linear regressions by expectation maximization. *IEEE Transactions on Information Theory*, 65(6):3515–3524, 2019.
- [LCZL20] Y. Li, Y. Chi, H. Zhang, and Y. Liang. Non-convex low-rank matrix recovery with arbitrary outliers via median-truncated gradient descent. *Information and Inference: A Journal of the IMA*, 9(2):289–325, 2020.
- [LL18] Y. Li and Y. Liang. Learning mixtures of linear regressions with nearly optimal complexity. In *Proceedings of the Conference On Learning Theory*, pages 1125–1144, 2018.
- [LMZ18] Y. Li, T. Ma, and H. Zhang. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In *Proceedings of the Conference On Learning Theory*, pages 2–47, 2018.
- [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.
- [MP20] A. Mazumdar and S. Pal. Recovery of sparse signals from a mixture of linear samples. *arXiv preprint arXiv:2006.16406*, 2020.
- [MPRP16] A. Maurer, M. Pontil, and B. Romera-Paredes. The benefit of multitask representation learning. *The Journal of Machine Learning Research*, 17(1):2853–2884, 2016.
- [MWCC20] 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. *Foundations of Computational Mathematics*, 20(3):451–632, 2020.
- [NNS<sup>+</sup>14] P. Netrapalli, U. Niranjan, S. Sanghavi, A. Anandkumar, and P. Jain. Non-convex robust pca. In *Advances in Neural Information Processing Systems*, pages 1107–1115, 2014.
- [PA18] D. Pimentel-Alarcón. Mixture matrix completion. In *Advances in Neural Information Processing Systems*, pages 2193–2203, 2018.
- [PKCS17] D. Park, A. Kyrillidis, C. Caramanis, and S. Sanghavi. Non-square matrix sensing without spurious local minima via the burer-monteiro approach. In *Artificial Intelligence and Statistics*, pages 65–74, 2017.
- [PL14] A. Pentina and C. Lampert. A pac-bayesian bound for lifelong learning. In *Proceedings of the International Conference on Machine Learning*, pages 991–999, 2014.
- [PLW<sup>+</sup>20] M. Peng, Y. Li, B. Wamsley, Y. Wei, and K. Roeder. cFIT: Integration and transfer learning of single cell transcriptomes, illustrated by fetal brain cell development. *bioRxiv*, 2020.
- [PY09] S. J. Pan and Q. Yang. A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(10):1345–1359, 2009.

- [QR78] R. E. Quandt and J. B. Ramsey. Estimating mixtures of normal distributions and switching regressions. *Journal of the American Statistical Association*, 73(364):730–738, 1978.
- [RFP10] B. Recht, M. Fazel, and P. A. Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Review*, 52(3):471–501, 2010.
- [Rou84] P. J. Rousseeuw. Least median of squares regression. *Journal of the American Statistical Association*, 79(388):871–880, 1984.
- [SBVDG10] N. Städler, P. Bühlmann, and S. Van De Geer. L1-penalization for mixture regression models. *Test*, 19(2):209–256, 2010.
- [Ser09] R. J. Serfling. *Approximation Theorems of Mathematical Statistics*, volume 162. John Wiley & Sons, 2009.
- [SJA16] H. Sedghi, M. Janzamin, and A. Anandkumar. Provable tensor methods for learning mixtures of generalized linear models. In *Artificial Intelligence and Statistics*, pages 1223–1231, 2016.
- [SL16] R. Sun and Z.-Q. Luo. Guaranteed matrix completion via non-convex factorization. *IEEE Transactions on Information Theory*, 62(11):6535–6579, 2016.
- [SLHH18] Y. Sun, Z. Liang, X. Huang, and Q. Huang. Joint map and symmetry synchronization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 251–264, 2018.
- [SQW18] J. Sun, Q. Qu, and J. Wright. A geometric analysis of phase retrieval. *Foundations of Computational Mathematics*, 18(5):1131–1198, 2018.
- [SS19a] Y. Shen and S. Sanghavi. Iterative least trimmed squares for mixed linear regression. In *Advances in Neural Information Processing Systems*, pages 6078–6088, 2019.
- [SS19b] Y. Shen and S. Sanghavi. Learning with bad training data via iterative trimmed loss minimization. In *Proceedings of the International Conference on Machine Learning*, pages 5739–5748, 2019.
- [SSZ17] J. Snell, K. Swersky, and R. Zemel. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems*, pages 4077–4087, 2017.
- [SWS20] V. Shah, X. Wu, and S. Sanghavi. Choosing the sample with lowest loss makes sgd robust. *arXiv preprint arXiv:2001.03316*, 2020.
- [TBS<sup>+</sup>16] S. Tu, R. Boczar, M. Simchowitz, M. Soltanolkotabi, and B. Recht. Low-rank solutions of linear matrix equations via procrustes flow. In *Proceedings of the International Conference on Machine Learning*, pages 964–973, 2016.
- [TJJ20] N. Tripuraneni, C. Jin, and M. I. Jordan. Provable meta-learning of linear representations. *arXiv preprint arXiv:2002.11684*, 2020.
- [TMC20] T. Tong, C. Ma, and Y. Chi. Accelerating ill-conditioned low-rank matrix estimation via scaled gradient descent. *arXiv preprint arXiv:2005.08898*, 2020.
- [Tur00] T. R. Turner. Estimating the propagation rate of a viral infection of potato plants via mixtures of regressions. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 49(3):371–384, 2000.
- [Ver18] R. Vershynin. *High-dimensional Probability: An Introduction with Applications in Data Science*, volume 47. Cambridge university press, 2018.
- [VT02] K. Viele and B. Tong. Modeling with mixtures of linear regressions. *Statistics and Computing*, 12(4):315–330, 2002.

- [Wed72] P.-Å. Wedin. Perturbation bounds in connection with singular value decomposition. *BIT Numerical Mathematics*, 12(1):99–111, 1972.
- [YC15] X. Yi and C. Caramanis. Regularized em algorithms: A unified framework and statistical guarantees. In *Advances in Neural Information Processing Systems*, pages 1567–1575, 2015.
- [YCS14] X. Yi, C. Caramanis, and S. Sanghavi. Alternating minimization for mixed linear regression. In *Proceedings of the International Conference on Machine Learning*, pages 613–621, 2014.
- [YCS16] X. Yi, C. Caramanis, and S. Sanghavi. Solving a mixture of many random linear equations by tensor decomposition and alternating minimization. *arXiv preprint arXiv:1608.05749*, 2016.
- [YL20] H. Yuan and Y. Liang. Learning entangled single-sample distributions via iterative trimming. volume 108 of *Proceedings of Machine Learning Research*, pages 2666–2676, Online, 26–28 Aug 2020. PMLR.
- [YPCR18] D. Yin, R. Pedarsani, Y. Chen, and K. Ramchandran. Learning mixtures of sparse linear regressions using sparse graph codes. *IEEE Transactions on Information Theory*, 65(3):1430–1451, 2018.
- [ZCL18] H. Zhang, Y. Chi, and Y. Liang. Median-truncated nonconvex approach for phase retrieval with outliers. *IEEE Transactions on Information Theory*, 64(11):7287–7310, 2018.
- [ZJD16] K. Zhong, P. Jain, and I. S. Dhillon. Mixed linear regression with multiple components. In *Advances in Neural Information Processing Systems*, pages 2190–2198, 2016.
- [ZL15] Q. Zheng and J. Lafferty. A convergent gradient descent algorithm for rank minimization and semidefinite programming from random linear measurements. In *Advances in Neural Information Processing Systems*, pages 109–117, 2015.
- [ZLW18] Z. Zhu, Q. Li, G. Tang, and M. B. Wakin. Global optimality in low-rank matrix optimization. *IEEE Transactions on Signal Processing*, 66(13):3614–3628, 2018.
- [ZQW20] Y. Zhang, Q. Qu, and J. Wright. From symmetry to geometry: Tractable nonconvex problems. *arXiv preprint arXiv:2007.06753*, 2020.
- [ZWYG18] X. Zhang, L. Wang, Y. Yu, and Q. Gu. A primal-dual analysis of global optimality in nonconvex low-rank matrix recovery. In *International conference on machine learning*, pages 5862–5871, 2018.