



Technical Report 109

# Solving a Mixture of Many Random Linear Equations by Tensor Decomposition and Alternating Minimization

**Research Supervisor:**  
Constantine Caramanis  
Wireless Networking and Communications Group

September 2016

Project title: *Models for High Dimensional Mixed Regression*

# Data-Supported Transportation Operations & Planning Center (D-STOP)

---

A Tier 1 USDOT University Transportation Center at The University of Texas at Austin



D-STOP is a collaborative initiative by researchers at the Center for Transportation Research and the Wireless Networking and Communications Group at The University of Texas at Austin.

## DISCLAIMER

The contents of this report reflect the views of the authors, who are responsible for the facts and the accuracy of the information presented herein. This document is disseminated under the sponsorship of the U.S. Department of Transportation's University Transportation Centers Program, in the interest of information exchange. The U.S. Government assumes no liability for the contents or use thereof.

**Technical Report Documentation Page**

1. Report No. <b>D-STOP/2016/109</b>		2. Government Accession No.		3. Recipient's Catalog No.	
4. Title and Subtitle <b>Solving a Mixture of Many Random Linear Equations by Tensor Decomposition and Alternating Minimization</b>				5. Report Date <b>September 2016</b>	
				6. Performing Organization Code	
7. Author(s) <b>Xinyang Yi, Constantine Caramanis, Sujay Sanghavi</b>				8. Performing Organization Report No. <b>Report 109</b>	
9. Performing Organization Name and Address <b>Data-Supported Transportation Operations &amp; Planning Center (D-STOP) The University of Texas at Austin 1616 Guadalupe Street, Suite 4.202 Austin, Texas 78701</b>				10. Work Unit No. (TRAIS)	
				11. Contract or Grant No. <b>DTRT13-G-UTC58</b>	
12. Sponsoring Agency Name and Address <b>Data-Supported Transportation Operations &amp; Planning Center (D-STOP) The University of Texas at Austin 1616 Guadalupe Street, Suite 4.202 Austin, Texas 78701</b>				13. Type of Report and Period Covered	
				14. Sponsoring Agency Code	
15. Supplementary Notes <b>Supported by a grant from the U.S. Department of Transportation, University Transportation Centers Program. Project Title: Models for High Dimensional Mixed Regression</b>					
16. Abstract <b>We consider the problem of solving mixed random linear equations with <math>k</math> components. This is the noiseless setting of mixed linear regression. The goal is to estimate multiple linear models from mixed samples in the case where the labels (which sample corresponds to which model) are not observed. We give a tractable algorithm for the mixed linear equation problem, and show that under some technical conditions, our algorithm is guaranteed to solve the problem exactly with sample complexity linear in the dimension, and polynomial in <math>k</math>, the number of components. Previous approaches have required either exponential dependence on <math>k</math>, or super-linear dependence on the dimension. The proposed algorithm is a combination of tensor decomposition and alternating minimization. Our analysis involves proving that the initialization provided by the tensor method allows alternating minimization, which is equivalent to EM in our setting, to converge to the global optimum at a linear rate.</b>					
17. Key Words <b>mixed random linear equations, mixed linear regression</b>			18. Distribution Statement <b>No restrictions. This document is available to the public through NTIS (<a href="http://www.ntis.gov">http://www.ntis.gov</a>): National Technical Information Service 5285 Port Royal Road Springfield, Virginia 22161</b>		
19. Security Classif.(of this report) <b>Unclassified</b>		20. Security Classif.(of this page) <b>Unclassified</b>		21. No. of Pages	22. Price

## Disclaimer

The contents of this report reflect the views of the authors, who are responsible for the facts and the accuracy of the information presented herein. Mention of trade names or commercial products does not constitute endorsement or recommendation for use.

## Acknowledgements

The authors recognize that support for this research was provided by a grant from the U.S. Department of Transportation, University Transportation Centers.

# Solving a Mixture of Many Random Linear Equations by Tensor Decomposition and Alternating Minimization

Xinyang Yi    Constantine Caramanis    Sujay Sanghavi

The University of Texas at Austin

{yixy,constantine}@utexas.edu    sanghavi@mail.utexas.edu

## Abstract

We consider the problem of solving mixed random linear equations with  $k$  components. This is the noiseless setting of mixed linear regression. The goal is to estimate multiple linear models from mixed samples in the case where the labels (which sample corresponds to which model) are not observed. We give a tractable algorithm for the mixed linear equation problem, and show that under some technical conditions, our algorithm is guaranteed to solve the problem exactly with sample complexity linear in the dimension, and polynomial in  $k$ , the number of components. Previous approaches have required either exponential dependence on  $k$ , or super-linear dependence on the dimension. The proposed algorithm is a combination of tensor decomposition and alternating minimization. Our analysis involves proving that the initialization provided by the tensor method allows alternating minimization, which is equivalent to EM in our setting, to converge to the global optimum at a linear rate.

## 1 Introduction

In this paper, we consider the following mixed linear equation problem. Suppose we are given  $n$  samples of response-covariate pairs  $\{(y_i, \mathbf{x}_i)\}_{i=1}^n$  that are determined by equations

$$y_i = \sum_{j=1}^k \langle \mathbf{x}_i, \boldsymbol{\beta}_j \rangle \mathbb{1}(z_i = j), \text{ for } i = 1, \dots, n, \quad (1)$$

where  $\mathbf{x}_i, \boldsymbol{\beta}_j \in \mathbb{R}^p$ ,  $\{\boldsymbol{\beta}_j\}$  are  $k$  model parameters corresponding to  $k$  different linear models, and  $z_i$  is the *unobserved* label of sample  $i$  indicating which model it is generated from. We assume random label assignment, i.e.,  $\{z_i\}$  are i.i.d. copies of a multinomial random variable  $Z$  that has distribution

$$\mathbb{P}[Z = j] = \omega_j, \text{ for } j = 1, 2, \dots, k. \quad (2)$$

Here  $\{\omega_j\}$  represent the weights of every linear model, and naturally satisfy  $\sum_{j \in [k]} \omega_j = 1$ . Our goal is to find parameters  $\{\boldsymbol{\beta}_j\}$  from mixed samples  $\{(y_i, \mathbf{x}_i)\}_{i=1}^n$ . While solving linear systems is straightforward, this problem, with the introduction of latent variables, is hard to solve in general. Work in [25] shows that the subset sum problem can be reduced to mixed linear equations in the case of  $k = 2$  and certain designs of  $\mathbf{x}_i$  and  $\boldsymbol{\beta}_j$ . Therefore, given  $\{(y_i, \mathbf{x}_i)\}_{i=1}^n$ , determining whether there exist two  $\boldsymbol{\beta}$ 's that satisfy (1) is NP-complete, and thus in general the  $k = 2$  case is already hard. In this paper, we consider the setting for general  $k$ , where the covariates  $\mathbf{x}_i$ 's are independently drawn from the standard Gaussian distribution:

$$\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p). \quad (3)$$

Under this random design, we provide a tractable algorithm for the mixed linear equation problem, and give sufficient conditions on the  $\beta_i$ 's under which we guarantee exact recovery with high probability.

The problem of solving mixed linear equations (or regression when each  $y_i$  is perturbed by a small amount of noise) arises in applications where samples are from a mixture of discriminative linear models and the interest is in parameter estimation. Mixed standard and generalized linear regression models are introduced in 1990s [23] and have become an important set of techniques for market segmentation [24]. These models have also been applied to study music perception [22] and health care demand [10]. See [11] for other related applications and datasets. Mixed linear regression is closely related to another classical model called hierarchical mixtures of experts [13], which also allows the distribution of labels to be adaptive to covariate vectors.

Due to the combinatorial nature of mixture models, popular approaches, including EM and gradient descent, are often based on local optimization and thus suffer from local minima. Indeed, to the best of our knowledge, there is no rigorous analysis of the convergence behavior of EM or other gradient descent-based methods for  $k \geq 3$ . Beyond real-world applications, the statistical limits of solving problem (1) by computationally efficient algorithms are even less well understood. This paper is motivated by this question: how many samples are necessary to recover  $\{\beta_j\}$  exactly and efficiently?

In a nutshell, we prove that under certain technical conditions, there exists an efficient algorithm for solving mixed linear equations with sample size  $\tilde{\mathcal{O}}(k^{10}p)$ , and we provide an algorithm which achieves this. Notably, the dependence on  $p$  is nearly linear and thus optimal up to some log factors. Our proposed algorithm has two phases. The first step is a spectral method called tensor decomposition, which is guaranteed to produce  $\varepsilon$ -close solutions with  $\mathcal{O}(1/\varepsilon^2)$  samples. In the second step, we apply an alternating minimization (AltMin) procedure to successively refine the estimation until exact recovery happens. As a key ingredient, we show that AltMin, as a non-convex optimization technique, enjoys linear convergence to the global optima when initialized closely enough to the optimal solution.

## 1.1 Comparison to Prior Art

The use of the method of moments for learning latent variable models can be dated back to Pearson's work [15] on estimating Gaussian mixtures. There is now an increasing interest in computing high order moments and leveraging tensor decomposition for parameter estimation in various mixture models including Hidden Markov Models [1], Gaussian mixtures [12], and topic models [3]. Following the same idea, we propose some novel moments for mixed linear equations, on which approximate estimation of parameters can be computed by tensor decomposition. Different from our moments, Chaganty and Liang [6] propose a method of regressing  $y_i^3$  against  $\mathbf{x}_i^{\otimes 3}$  to estimate a certain third-moment tensor of mixed linear regression under bounded and random covariates. Because of performing regression in the lifted space with dimension  $p^3$ , their method suffers from much higher sample complexity  $\mathcal{O}(p^6)$  compared to our results, while the latter builds on a different covariate assumption (3).

Mixed linear equation/regression with two components is now well understood. In particular, our earlier work [25] proves the local convergence of AltMin for mixed linear equations with two components. Through a convex optimization formula, work in [8] establishes the minimax optimal statistical rates under stochastic and deterministic noises. Notably, Balakrishnan et al. [4] develop a framework for analyzing the local convergence of expectation-maximization (EM), i.e., EM is

guaranteed to produce statistically consistent points with good initializations. In the case of mixed linear regression with  $\beta_1 = -\beta_2$  and Gaussian noise with variance  $\sigma^2$ , applying the framework leads to estimation error  $\tilde{\mathcal{O}}(\sqrt{(\sigma^2 + \|\beta_1\|_2)p/n})$ . Even in the case of no noise ( $\sigma = 0$ ), their results do not imply exact recovery. Moreover, it is unclear how to apply the framework to the case of  $k \geq 3$  components. It is obvious that AltMin is equivalent to EM in the noiseless setting. Our analysis of AltMin takes a step further towards understanding EM in the case of multiple latent clusters.

Beyond linear models, learning mixture of generalized linear models is recently studied in [19] and [16]. Specifically, [19] proposes a spectral method based on second order moments for estimating the subspace spanned by model parameters. Later on, Sedghi et al. [16] construct specific third order moments that allow tensor decomposition to be applied to estimate individual vectors. In detail, when  $k = \mathcal{O}(1)$ , they show that obtaining recovery error  $\varepsilon$  requires sample size  $n = \tilde{\mathcal{O}}(p^3/\varepsilon^2)$ . In a more recent update [17] of their paper, they establish the same sample complexity for mixed linear regression using different moments, which we realize coincide with ours during the preparation of this paper. Nevertheless, we perform a sharper analysis that leads to a near-linear-in- $p$  sample complexity  $n = \tilde{\mathcal{O}}(p/\varepsilon^2)$ .

Conceptually, we establish the power of combining spectral method and likelihood based estimation for learning latent variable models. Spectral method excludes most bad local optima on the surface of likelihood loss, and as a consequence, it becomes much easier for non-convex local search methods such as EM and AltMin, to find statistically efficient solutions. Such phenomenon in the context of mixed linear regression is observed empirically in [6]. We provide a theoretical explanation in this paper. It is worth mentioning the applications of such idea in other problems including crowdsourcing [26], phase retrieval (e.g. [5, 9]) and matrix completion (e.g. [14, 18, 7]). Most of these works focus on estimating bilinear or low rank structures. In the context of crowdsourcing, work in [26] shows that performing one step of EM can achieve optimal rate given good initialization. In contrast, we establish an explicit convergence trajectory of multiple steps of AltMin for our problem. It would be interesting to study the convergence path of AltMin or EM for other latent variable models.

## 1.2 Notation and Outline

We lay down some notations commonly used throughout this paper. For counting number  $k$ , we use  $[k]$  to denote the set  $\{1, 2, \dots, k\}$ . We let  $a \vee b$ ,  $a \wedge b$  denote  $\max\{a, b\}$ ,  $\min\{a, b\}$  respectively. For sub-Gaussian random variable  $X$ , we denote its  $\psi_2$ -Orlicz norm [20] by  $\|X\|_{\psi_2}$ , i.e.,

$$\|X\|_{\psi_2} := \inf \{z \in (0, \infty) \mid \mathbb{E}[\psi_2(|X|/z)] \leq 1\},$$

where  $\psi_2(x) = \exp(x^2) - 1$ . For vector  $\mathbf{a} \in \mathbb{R}^p$ , we use  $\|\mathbf{a}\|_q$  to denote the standard  $\ell_q$  norm of  $\mathbf{a}$ . For matrix  $\mathbf{A} \in \mathbb{R}^{p_1 \times p_2}$ , we use  $\sigma_k(\mathbf{A})$  to denote its  $k$ -th largest singular value. We also commonly use  $\sigma_{\max}(\mathbf{A})$ ,  $\sigma_{\min}(\mathbf{A})$  to denote  $\sigma_1(\mathbf{A})$  and  $\sigma_{p_1 \wedge p_2}(\mathbf{A})$ . In particular, we denote the operator norm of matrix  $\mathbf{A}$  as  $\|\mathbf{A}\|_{op}$ . We also use  $\|\mathbf{T}\|_{op}$  to denote the operator norm of symmetric third order tensor  $\mathbf{T} \in \mathbb{R}^{p \times p \times p}$ , namely

$$\|\mathbf{T}\|_{op} := \sup_{\mathbf{u} \in \mathbb{S}^{p-1}} |\mathbf{T}(\mathbf{u}, \mathbf{u}, \mathbf{u})|.$$

Here,  $\mathbf{T}(\mathbf{A}, \mathbf{B}, \mathbf{C})$  denotes the multi-linear matrix multiplication of  $\mathbf{T}$  by  $\mathbf{A} \in \mathbb{R}^{p \times p_1}$ ,  $\mathbf{B} \in \mathbb{R}^{p \times p_2}$ ,  $\mathbf{C} \in \mathbb{R}^{p \times p_3}$ , namely,

$$(\mathbf{T}(\mathbf{A}, \mathbf{B}, \mathbf{C}))_{(m,n,t)} = \sum_{i,j,k \in [p]} \mathbf{T}_{(i,j,k)} \mathbf{A}_{(i,m)} \mathbf{B}_{(j,n)} \mathbf{C}_{(k,t)}, \quad \text{for all } (m, n, t) \in [p_1] \times [p_2] \times [p_3].$$

For two sequences  $f(n), g(n)$  indexed by  $n \in \mathbb{N}$ , we write  $f(n) = \mathcal{O}(g(n))$  to mean there exists a constant  $C > 0$  such that  $f(n) \leq Cg(n)$  for all  $n \in \mathbb{N}$ . By  $f(n) = \tilde{\mathcal{O}}(g(n))$ , we mean there exist constants  $C, C' > 0$  such that  $f(n) \leq Cg(n) \cdot (\log n)^{C'}$ . We also use  $f(n) \lesssim g(n)$  as shorthand for  $f(n) = \mathcal{O}(g(n))$ . Similarly, we say  $f(n) \gtrsim g(n)$  if  $g(n) = \mathcal{O}(f(n))$ .

The rest of this paper is organized as follows. In Section 2, we describe the specific details of our two-phase algorithm for solving mixed linear equations. We present the theoretical results of initialization and AltMin in Section 3.1 and 3.2 respectively. We combine these two parts and give the overall sample and time complexities for exact recovery in Section 3.3. We provide the experimental results in Section 4. All proofs are collected in Section 5.

## 2 Algorithm

A natural idea to solve problem (1) is to apply an alternating minimization (AltMin) procedure between parameters  $\{\beta_j\}$  and labels  $\{z_i\}$ : (1) Given  $\{\beta_j\}$ , assign the labels for each sample by choosing a model  $\beta$  that has minimal recovery error  $|y_i - \langle \mathbf{x}_i, \beta \rangle|$ ; (2) When labels are available, each parameter is updated by applying the method of least square optimization to samples with the corresponding labels. One can show that in our setting, alternating minimization is equivalent to Expectation-Maximization (EM), which is one of the most important algorithms for inference in latent variable models. In general, similar to EM, AltMin is vulnerable to local optima. Our experiment (see Figure 1) demonstrates that even under random setting  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ , AltMin with random initializations fails to exactly recover each  $\beta_j$  with significantly large probability.

To overcome the local-optima issue of AltMin, our algorithm consists of two stages. The first stage builds on carefully designed moments of samples, and aims to find rough estimates of  $\{\beta_j\}$ . Starting with the initialization, the second stage involves using AltMin to successively refine the estimates. In the following, we describe these two steps with more details.

### 2.1 Tensor Decomposition

In the first step, we use method of moments to compute initial estimates of  $\{\beta_j\}$ . Consider moments  $m_0 \in \mathbb{R}, \mathbf{m}_1 \in \mathbb{R}^p, \mathbf{M}_2 \in \mathbb{R}^{p \times p}$  and  $\mathbf{M}_3 \in \mathbb{R}^{p \times p \times p}$  as

$$m_0 := \frac{1}{n} \sum_{i=1}^n y_i^2, \quad \mathbf{m}_1 := \frac{1}{6n} \sum_{i=1}^n y_i^3 \mathbf{x}_i, \quad (4)$$

$$\mathbf{M}_2 := \frac{1}{2n} \sum_{i=1}^n y_i^2 \mathbf{x}_i \otimes \mathbf{x}_i - \frac{1}{2} m_0 \cdot \mathbf{I}_p, \quad (5)$$

$$\mathbf{M}_3 := \frac{1}{6n} \sum_{i=1}^n y_i^3 \mathbf{x}_i \otimes \mathbf{x}_i \otimes \mathbf{x}_i - \mathcal{T}(\mathbf{m}_1), \quad (6)$$

where  $\mathcal{T}(\cdot)$  is a mapping from  $\mathbb{R}^p$  to  $\mathbb{R}^{p \times p \times p}$  with form

$$\mathcal{T}(\mathbf{m}_1) := \sum_{i \in [p]} \mathbf{m}_1 \otimes \mathbf{e}_i \otimes \mathbf{e}_i + \mathbf{e}_i \otimes \mathbf{m}_1 \otimes \mathbf{e}_i + \mathbf{e}_i \otimes \mathbf{e}_i \otimes \mathbf{m}_1.$$

It is reasonable to choose these moments because of the next result, which shows that the expectations of  $M_2$  and  $M_3$  contain the structure of  $\{\beta_j\}$ . See Section 5.1 for its proof.



**Lemma 1** (Moment Expectation). *Consider the random model for mixed linear equations given in (1), (2) and (3). For moments  $\mathbf{M}_2$  and  $\mathbf{M}_3$  in (5) and (6), we have*

$$\mathbb{E}[\mathbf{M}_2] = \sum_{j=1}^k \omega_j \cdot \boldsymbol{\beta}_j \otimes \boldsymbol{\beta}_j, \quad (7)$$

$$\mathbb{E}[\mathbf{M}_3] = \sum_{j=1}^k \omega_j \cdot \boldsymbol{\beta}_j \otimes \boldsymbol{\beta}_j \otimes \boldsymbol{\beta}_j. \quad (8)$$

With the special structure given on the right hand sides of (7) and (8), tensor decomposition techniques can discover  $\{(\omega_j, \boldsymbol{\beta}_j)\}$  in three steps under a *non-degeneracy condition* (see Condition 1). First, apply SVD on  $\mathbb{E}[\mathbf{M}_2]$  to compute a whitening matrix  $\mathbf{W} \in \mathbb{R}^{p \times k}$  such that  $\mathbf{W}^\top \mathbb{E}[\mathbf{M}_2] \mathbf{W} = \mathbf{I}_p$ . Then we use  $\mathbf{W}$  to transform  $\mathbb{E}[\mathbf{M}_3]$  into an orthogonal tensor  $\mathbb{E}[\mathbf{M}_3](\mathbf{W}, \mathbf{W}, \mathbf{W})$ , which is further decomposed into eigenvalue/eigenvector pairs by robust tensor power method (Algorithm 2). Lastly,  $\{(\omega_j, \boldsymbol{\beta}_j)\}$  can be reconstructed by applying simple linear transformation upon the previously discovered spectral components from  $\mathbb{E}[\mathbf{M}_3](\mathbf{W}, \mathbf{W}, \mathbf{W})$ . With sufficient amount of samples, it is reasonable to believe that  $\mathbf{M}_2$  and  $\mathbf{M}_3$  are close to their expectations such that the stability of tensor decomposition will lead to good enough estimates. For the ease of analysis, we need to ensure the independence between whitening matrix  $\mathbf{W}$  and  $\mathbf{M}_3$ . Accordingly, we split the samples used in initialization into two disjoint parts for computing  $\{m_0, \mathbf{M}_2\}$  and  $\{\mathbf{m}_1, \mathbf{M}_3\}$  respectively. We present the details in Algorithm 1.

---

**Algorithm 1** Initialization via Tensor Factorization

---

**INPUT:** Samples  $\{(y_i, \mathbf{x}_i)\}_{i=1}^n$ .

- 1: Randomly split samples into two disjoint parts  $\{(y_i, \mathbf{x}_i)\}_{i=1}^{n_1}$  and  $\{(y'_i, \mathbf{x}'_i)\}_{i=1}^{n_2}$ .
- 2:  $m_0 \leftarrow \frac{1}{n_1} \sum_{i=1}^{n_1} y_i^2$ ,  $\mathbf{m}_1 \leftarrow \frac{1}{6n_2} \sum_{i=1}^{n_2} y_i'^3 \mathbf{x}'_i$ .
- 3:  $\mathbf{M}_2 \leftarrow \frac{1}{2n_1} \sum_{i=1}^{n_1} y_i^2 \mathbf{x}_i \otimes \mathbf{x}_i - \frac{1}{2} m_0 \cdot \mathbf{I}_p$ ,  $\mathbf{M}_3 \leftarrow \frac{1}{6n_2} \sum_{i=1}^{n_2} y_i'^3 \mathbf{x}'_i \otimes \mathbf{x}'_i \otimes \mathbf{x}'_i - \mathcal{T}(\mathbf{m}_1)$ .
- 4: Compute an SVD of the best rank  $k$  approximation of  $\mathbf{M}_2$  as  $\mathbf{U}\boldsymbol{\Sigma}\mathbf{U}^\top$ , where  $\mathbf{U} \in \mathbb{R}^{p \times k}$ . Compute whitening matrix  $\mathbf{W} \leftarrow \mathbf{U}\boldsymbol{\Sigma}^{-1/2}$ .
- 5:  $\widetilde{\mathbf{M}}_3 \leftarrow \mathbf{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})$ .
- 6: Run robust tensor power method (Algorithm 2) on  $\widetilde{\mathbf{M}}_3$  to obtain  $k$  eigenvalue/eigenvector pairs  $\{(\widetilde{\omega}_j, \widetilde{\boldsymbol{\beta}}_j)\}_{j=1}^k$ .
- 7:  $\omega_j^{(0)} \leftarrow 1/\widetilde{\omega}_j^2$ ,  $\boldsymbol{\beta}_j^{(0)} \leftarrow \widetilde{\omega}_j (\mathbf{W}^\top)^\dagger \widetilde{\boldsymbol{\beta}}_j$ , for all  $j \in [k]$ .<sup>1</sup>

**OUTPUT:**  $\{(\omega_j^{(0)}, \boldsymbol{\beta}_j^{(0)})\}_{j=1}^k$ .

---

<sup>1</sup> $(\mathbf{W}^\top)^\dagger$  denotes the Moore-Penrose pseudoinverse of  $\mathbf{W}^\top$ , i.e.,  $\mathbf{W}(\mathbf{W}^\top \mathbf{W})^{-1}$ .

---

**Algorithm 2** Robust Tensor Power Method (Algorithm 1 in [2])
 

---

**INPUT:** Symmetric tensor  $\mathbf{T} \in \mathbb{R}^{k \times k \times k}$ . Parameters  $L, N$ .

- 1: **for**  $j = 1, \dots, k$  **do**
- 2:   **for**  $l = 1, \dots, L$  **do**
- 3:     Draw  $\beta_0^{(l)}$  uniformly at random from  $\mathbb{S}^{k-1}$ .
- 4:     **for**  $t = 0, \dots, N - 1$  **do**

$$\beta_{t+1}^{(l)} \leftarrow \mathbf{T}(\mathbf{I}_k, \beta_t^{(l)}, \beta_t^{(l)}), \quad \beta_{t+1}^{(l)} \leftarrow \beta_{t+1}^{(l)} / \left\| \beta_{t+1}^{(l)} \right\|_2. \quad (9)$$

- 5:     **end for**
- 6:   **end for**
- 7:    $l^* \leftarrow \arg \max_{l \in [L]} \mathbf{T}(\beta_N^{(l)}, \beta_N^{(l)}, \beta_N^{(l)})$ .
- 8:   Do  $N$  power updates (9) starting from  $\beta_N^{(l^*)}$  to obtain  $\tilde{\beta}_j$ . Let  $\tilde{\omega}_j \leftarrow \mathbf{T}(\tilde{\beta}_j, \tilde{\beta}_j, \tilde{\beta}_j)$ .
- 9:    $T \leftarrow T - \tilde{\omega}_j \tilde{\beta}_j^{\otimes 3}$ .
- 10: **end for**

**OUTPUT:**  $\{(\tilde{\omega}_j, \tilde{\beta}_j)\}_{j=1}^k$ .

---

## 2.2 Alternating Minimization

The motivation for using AltMin is to consider the least-square loss function below

$$\mathcal{L}_n(\{\beta_j\}) := \min_{z_1, \dots, z_n \in [k]} \sum_{i=1}^n \sum_{j=1}^k (y_i - \langle \mathbf{x}_i, \beta_j \rangle)^2 \mathbf{1}(z_i = j).$$

The minimization over discrete labels  $\{z_i\}$  makes the above loss function non-convex and yields hardness of solving mixed linear equations in general. A natural idea to minimize  $\mathcal{L}_n$  is by minimizing  $\{z_i\}$  and  $\{\beta_j\}$  alternatively and iteratively. Given initial estimates  $\{\beta_j^{(0)}\}$ , each iteration  $t = 0, 1, \dots$  consists of the following two steps:

- **Label Assignment:** Pick the model that has the smallest reconstruction error for each sample

$$z_i^{(t)} = \arg \min_{j \in [k]} |y_i - \langle \mathbf{x}_i, \beta_j^{(t)} \rangle|. \quad (10)$$

- **Parameter Update:**

$$\beta_j^{(t+1)} = \arg \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n (y_i - \langle \mathbf{x}_i, \beta \rangle)^2 \mathbf{1}(z_i^{(t)} = j). \quad (11)$$

AltMin runs quickly and is thus favored in practice. However, as we discussed before, its convergence to global optima is commonly intractable. In order to alleviate such issue, we already discussed how to construct good initial estimates by method of moments. Here, we introduce another ingredient—*resampling*—for making the analysis of AltMin tractable. The key idea is to split all samples into multiple disjoint subsets and use a fresh piece of samples in each iteration. While slightly inefficient regarding sample complexity, this trick decouples the probabilistic dependence between

two successive estimates  $\{\beta_j^{(t)}\}$  and  $\{\beta_j^{(t+1)}\}$ , and thus makes our analysis hold. The details are presented in Algorithm 3.

---

**Algorithm 3** Alternating Minimization with Resampling

---

**INPUT:** Samples  $\{(y_i, \mathbf{x}_i)\}_{i=1}^n$ , initial estimates  $\{\beta_j^{(0)}\}$ , number of iterations  $T$ .

1: Split all samples into  $T$  disjoint subsets  $\{(y_i^{(t)}, \mathbf{x}_i^{(t)})\}_{i=1}^{n/T}, t = 0, 1, \dots, T - 1$ , with equal size.

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

3:

$$z_i^{(t)} \leftarrow \arg \min_{j \in [k]} |y_i^{(t)} - \langle \mathbf{x}_i^{(t)}, \beta_j^{(t)} \rangle|, \text{ for all } i \in [n].$$

4:

$$\beta_j^{(t+1)} \leftarrow \arg \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^{n/T} (y_i^{(t)} - \langle \mathbf{x}_i^{(t)}, \beta \rangle)^2 \mathbb{1}(z_i^{(t)} = j), \text{ for all } j \in [k].$$

5: **end for**

**OUTPUT:**  $\{\beta_j^{(T)}\}_{j=1}^k$ .

---

### 3 Theoretical Results

In this section, we provide the theoretical guarantees of Algorithm 1 and 3. For simplicity, we assume the  $\ell_2$  norm of  $\beta_j$  is at most 1, i.e.,

$$\max_{j \in [k]} \|\beta_j\|_2 = 1.$$

Moreover, we impose the following non-degeneracy condition on  $\{\beta_j\}$ .

**Condition 1** (Non-degeneracy). *Parameters  $\beta_1, \dots, \beta_k$  are linearly independent and all weights  $\omega_j$  are strictly greater than 0, namely*

$$\underline{\omega} := \min_{j \in [k]} \omega_j > 0.$$

Under the above condition,  $\overline{\mathbf{M}}_2 = \sum_{j \in [k]} \omega_j \beta_j \otimes \beta_j$  has rank  $k$ , which leads to

$$\sigma_k := \sigma_k(\overline{\mathbf{M}}_2) > 0.$$

We use  $\Delta$  to denote the minimum distance between any two parameters, namely

$$\Delta := \min_{i, j \in [k], i \neq j} \|\beta_i - \beta_j\|_2.$$

The above three quantities  $(\underline{\omega}, \sigma_k, \Delta)$  represent the hardness of our problem, and will appear in the results of our analysis. For estimates  $\{\hat{\beta}_j\}$ , we define the estimation error  $\mathcal{E}(\{\hat{\beta}_j\})$  as

$$\mathcal{E}(\{\hat{\beta}_j\}) := \inf_{\pi} \sup_{j \in [k]} \left\| \hat{\beta}_j - \beta_{\pi(j)} \right\|_2, \quad (12)$$

where the infimum is taken over all permutations  $\pi(\cdot)$  on  $[k]$ .

### 3.1 Analysis of Tensor Decomposition

Our first result, proved in Section 5.4, provides a guarantee of Algorithm 1.

**Theorem 1** (Tensor Decomposition). *Consider Algorithm 1 for initial estimation of  $\{\beta_j\}$ . Pick any  $\delta \in (0, 1)$ . There exist constants  $C_i$  such that the following holds. Pick any  $\varepsilon \in (0, C_1/k)$ . If*

$$n_1 \geq C_2 \left( \frac{p \log(12k/\delta) \log^2 n_1}{\underline{\omega} \sigma_k^5 \varepsilon^2} \vee \frac{k}{\underline{\omega} \delta} \right) \quad \text{and} \quad n_2 \geq C_3 \left( \frac{(k^2 \vee p) \log(12k/\delta) \log^3 n_2}{\underline{\omega} \sigma_k^3 \varepsilon^2} \vee \frac{k}{\underline{\omega} \delta} \right), \quad (13)$$

then with probability at least  $1 - \delta$ , the output  $\{\beta_j^{(0)}\}$  satisfy

$$\mathcal{E}(\{\beta_j^{(0)}\}) \leq \varepsilon.$$

Theorem 1 shows that  $n_1, n_2$  have inverse dependencies on  $\underline{\omega}, \sigma_k$ . In the well balanced setting, we have  $\underline{\omega} = \Omega(1/k)$ . In general,  $\sigma_k$  can be quite small, especially in the case where some parameter  $\beta$  almost lies in the subspace spanned by the rest  $k - 1$  parameters and has a very small magnitude along the orthogonal direction. Below we provide a sufficient condition under which  $\sigma_k$  has a well established lower bound.

**Condition 2** (Nearly Orthonormal Condition( $\eta, \gamma$ )). *For all  $j \in [k]$ ,  $\|\beta_j\|_2 \geq 1 - \eta$ . Moreover,  $|\langle \beta_i, \beta_j \rangle| \leq \gamma$  for all  $i, j \in [k], i \neq j$ .*

Under the above condition, the next result provides a lower bound of  $\sigma_k$ . See Section 5.2 for the proof.

**Lemma 2.** *Suppose  $\{\beta_j\}$  satisfy the nearly orthonormal condition with  $\eta, \gamma$ . Then we have*

$$\sigma_k \geq \underline{\omega}(1 - \eta - k\gamma).$$

In the following discussion, we focus on balanced clusters, i.e.,  $\underline{\omega} \gtrsim 1/k$ . We also assume that  $\{\beta_j\}$  satisfy Condition 2 with  $\eta \lesssim 1$  and  $\gamma \lesssim 1/k$ , which leads to  $\sigma_k = \Omega(\underline{\omega})$  according to Lemma 2. Now we provide two remarks for Theorem 1.

*Remark 1* (Sample Complexity). We treat  $\delta$  in Theorem 1 as a constant. Then (13) implies that  $n = n_1 + n_2 = \mathcal{O}(\varepsilon^{-2} k^6 p \log k \log^3(p/\varepsilon))$  is sufficient to guarantee that the estimates produced by Algorithm 1 have accuracy at most  $\varepsilon$ . Moreover, we have  $n_1 = \tilde{\mathcal{O}}(\varepsilon^{-2} k^6 p)$ ,  $n_2 = \tilde{\mathcal{O}}(\varepsilon^{-2}(k^6 + k^4 p))$ , which indicates that more samples are required to compute  $\mathbf{M}_2$  than  $\mathbf{M}_3$ . To provide some intuitions why this conclusion makes sense, note that the estimation accuracy of  $\bar{\mathbf{M}}_2$  determines the accuracy of identifying the subspace spanned by  $\{\beta_j\}$  in the original  $p$ -dimensional space. While  $\mathbf{M}_3$  has higher order, it is only required to concentrate well on a  $k$ -dimensional subspace computed from  $\mathbf{M}_2$  thanks to the whitening procedure. It turns out subspace accuracy has a more critical impact on the final error and needs to be sharpened with more samples.

*Remark 2* (Time Complexity). Except the line 6 in Algorithm 1, the other steps have total complexity  $\mathcal{O}(n(p^2 + k^3))$ . Note that it's not necessary to compute  $\mathbf{M}_3$  directly since we can compute  $\tilde{\mathbf{M}}_3$  from whitened covariate vectors  $\mathbf{W}^\top \mathbf{x}_i$ . Running time of robust tensor power method is  $\mathcal{O}(k^4 NL)$ . According to Lemma 4, it is sufficient to set  $N = \mathcal{O}(\log k + \log \log(1/\varepsilon))$  and  $L = \mathcal{O}(\text{poly}(k))$  for some polynomial function  $\text{poly}(\cdot)$ . When  $k$  is large enough,  $L$  can be very close to be linear in  $k$  (see Theorem 5.1 in [2] for details). Roughly, we take  $L = \mathcal{O}(k^2)$ , which gives the running time of Algorithm 2 as  $\mathcal{O}(k^6 \log k)$  when  $\varepsilon \gtrsim \text{poly}(1/k)$ . Therefore, the overall complexity of Algorithm 1 is  $\mathcal{O}(n(p^2 + k^3) + k^6 \log k)$ .

### 3.2 Analysis of Alternating Minimization

Now we turn to the analysis of Algorithm 3. Let  $\varepsilon_0 := \mathcal{E}(\{\beta_j^{(0)}\})$ .

**Theorem 2** (Alternating Minimization). *Consider Algorithm 3 for successively refining estimation of  $\{\beta_j\}$ . Pick any  $\delta \in (0, 1)$ . There exist constants  $C_i$  such that the following holds. Suppose*

$$\varepsilon_0 \leq C_1 \left( \frac{1}{k^2} \wedge \underline{\omega} \right) \Delta, \quad p \geq \log(2k^2T/\delta),$$

and  $n$  satisfies

$$n/T \geq C_2 \left( \frac{kp}{\underline{\omega}} \vee \frac{\log(8k^2T/\delta)}{\underline{\omega}^2} \right). \quad (14)$$

With probability at least  $1 - \delta$ ,  $\{\beta_j^{(t)}\}$  satisfies

$$\mathcal{E}(\{\beta_j^{(t)}\}) \leq \left( \frac{1}{2} \right)^t \cdot \varepsilon_0, \quad \text{for } t = 1, \dots, T.$$

See Section 5.5 for the proof of the above result. Theorem 3 suggests that with good enough initialization, iterates  $\{\beta_j^{(t)}\}$  have at least linear convergence to the ground truth parameters. Due to the fast convergence, it is sufficient to set  $T = \mathcal{O}(\log(1/\varepsilon))$  to obtain estimation with accuracy  $\varepsilon$ . In the case of well balanced clusters, i.e.  $\underline{\omega} \gtrsim 1/k$ ,  $\varepsilon_0$  is required to be  $\mathcal{O}(\Delta/k^2)$  in order to guarantee the convergence to global optima. Next, we give two remarks for sample and time complexities. In our discussion, we assume  $\underline{\omega} \gtrsim 1/k$  and that  $\delta$  is a small constant.

*Remark 3* (Sample Complexity). For accuracy  $\varepsilon$ , it is sufficient to have  $n = \mathcal{O}(k^2p \log(1/\varepsilon))$  when  $p$  satisfies  $p \gtrsim \log k + \log \log(1/\varepsilon)$ . Compared to the sample complexity of tensor decomposition, AltMin avoids the high-order polynomial factor of  $k$ . Moreover, it also changes the dependence on  $\varepsilon$  from  $1/\varepsilon^2$  to  $\log(1/\varepsilon)$ , which is a big save especially when we focus on exact recovery, which can happen as we show in the next section, after one step of AltMin when  $\varepsilon \lesssim 1/p$ . Notably, the statistical efficiency comes from a good initialization provided by tensor/spectral method. On one hand, AltMin alleviates the statistical inefficiency of spectral method; on the other hand, spectral method resolves the algorithmic intractability of AltMin.

*Remark 4* (Time Complexity). Each iteration of AltMin has time complexity  $\mathcal{O}(np^2/T + kp^3)$ . Hence, the overall running time is  $\mathcal{O}(np^2 + kp^3 \log(1/\varepsilon))^2$ . Using the minimum requirement of  $n$ , we obtain complexity  $\tilde{\mathcal{O}}(k^2p^3)$ . Recall that solving linear regression by most practical algorithms has complexity  $\mathcal{O}(p^3)$ . Therefore, even labels are available, solving  $k$  sets of linear equations requires time  $\mathcal{O}(kp^3)$ . AltMin almost has an extra factor  $k$  as the price for addressing latent variables.

### 3.3 Exact Recovery and Overall Guarantee

We now consider putting the previous analysis of tensor decomposition and AltMin together to show exact recovery of  $\{\beta_j\}$ .

---

<sup>2</sup>Factor  $p^3$  in the second term stands for the complexity of inverting a  $p$ -by- $p$  matrix by Gauss-Jordan elimination. It can be further reduced by more complicated algorithms such as Strassen algorithm that has  $\mathcal{O}(p^{2.807})$ .

**Lemma 3.** Pick any  $\delta \in (0, 1)$ . For any fixed estimates  $\{\widehat{\beta}_j\}_{j=1}^k$  and some constant  $C$ , if

$$n \geq C \frac{1}{\underline{\omega}} (p \vee \log(k/\delta)) \quad \text{and} \quad \mathcal{E}(\{\widehat{\beta}_j\}) \leq \frac{\delta}{4nk} \Delta,$$

Running one step of alternating minimization according to (10) and (11) using  $n$  samples and initial guess  $\{\widehat{\beta}_j\}$  produces true parameters  $\{\beta_j\}$  with probability at least  $1 - \delta$ .

We provide the proof of the above result in Section 5.6. Putting all ingredients together, we have the following overall guarantee:

**Corollary 1** (Exact Recovery). Consider splitting  $n$  samples from (1) into two disjoint sets with size  $n_{\text{init}}, n_{\text{alt}}$  as inputs of Algorithm 1 and 3 for solving mixed linear equations as a two-stage method. Pick any  $\delta \in (0, 1)$ . There exist constants  $C_i$  such that the following holds. If we choose  $T = C_1 \log(kn_{\text{alt}}/\delta)$  in Algorithm 3, and  $(n_{\text{init}}, n_{\text{alt}}, p)$  satisfy

$$n_{\text{init}} \geq C_2 \left( \frac{(k^4 + 1/\underline{\omega}^2)(p/\sigma_k^2 + k^2 + p) \log(k/\delta)}{\underline{\omega} \sigma_k^3 \Delta^2} \log^3(n_{\text{init}}) + \frac{k}{\underline{\omega} \delta} \right),$$

$$n_{\text{alt}} \geq C_3 \left( \frac{kp}{\underline{\omega}} + \frac{p}{\underline{\omega}^2} \right) \log(kn_{\text{alt}}/\delta),$$

and

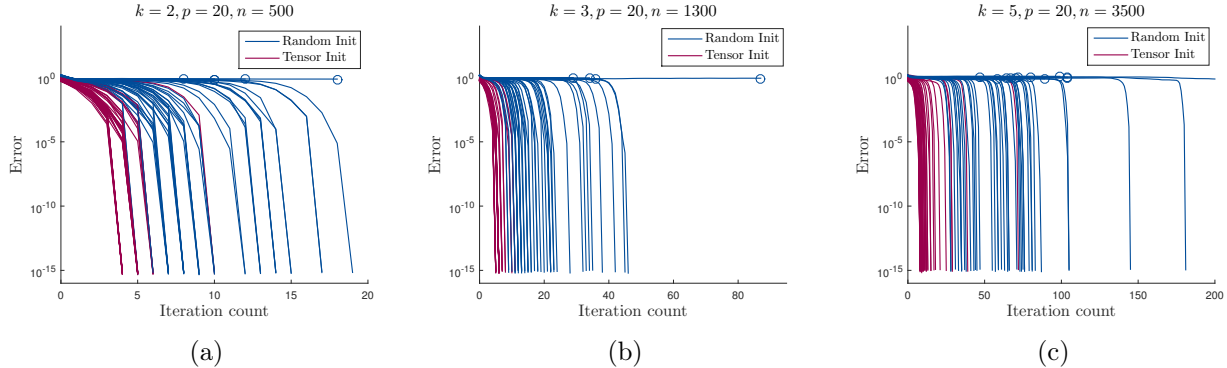
$$p \geq C_4 \left[ \log \left( \frac{k}{\delta} \right) + \log \log \left( \frac{kn_{\text{alt}}}{\delta} \right) \right],$$

then with probability at least  $1 - \delta$ , we have exact recovery, i.e.  $\{\beta_j^{(T)}\}_{j=1}^k = \{\beta_j\}_{j=1}^k$ .

The proof is provided in Section 5.3. When  $\underline{\omega} \gtrsim 1/k$  and Condition 2 holds with  $\gamma \lesssim 1$  and  $\eta \lesssim 1/k$  ( $\Delta \gtrsim 1$  in the case), Corollary 1 implies that  $n = n_{\text{init}} + n_{\text{alt}} = \mathcal{O}(k^{10} p \log k \log^3 p)$  is enough for exact recovery with high probability, say 99%. With this amount of samples, Remarks 2 and 4 give the overall time complexity as  $\mathcal{O}(k^{10} p (p^2 + k^3) \log k \log^3 p)$ . Note that solving  $k$  sets of linear equations (labels are known) needs at least  $kp$  samples, and usually requires time  $\mathcal{O}(kp^3)$ . Hence, under the aforementioned setting, our two-stage algorithm is nearly optimal in  $p$  with respect to sample and time complexities.

## 4 Numerical Results

In this section, we provide some numerical results to demonstrate the empirical performance of the proposed method (combination of Algorithms 1 and 3) for solving mixed linear equations, and also compare it with random initialized Alternating minimization (AltMin). All algorithms are implemented in MATLAB. While sample-splitting is useful for our theoretical analysis, we find it unnecessary in practice. Therefore, we remove the sample-splittings in Algorithms 1 and 3, and use the whole sample set in the entire process. AltMin is implemented to terminate when the label assignment no longer changes or the maximal number of iterations  $T$  is reached. In all experiments, we set  $T = 200$ .



**Figure 1.** Plot of estimation error (log scale) versus number of iterations in AltMin. Each panel shows 50 trials for random and tensor initializations respectively. The circle markers indicate the terminations of AltMin due to local minima, i.e., the label assignments do not change in two consecutive iterations. Tensor decomposition is implemented with  $L = 200k^2$ ,  $N = 20 \log(k)$ .

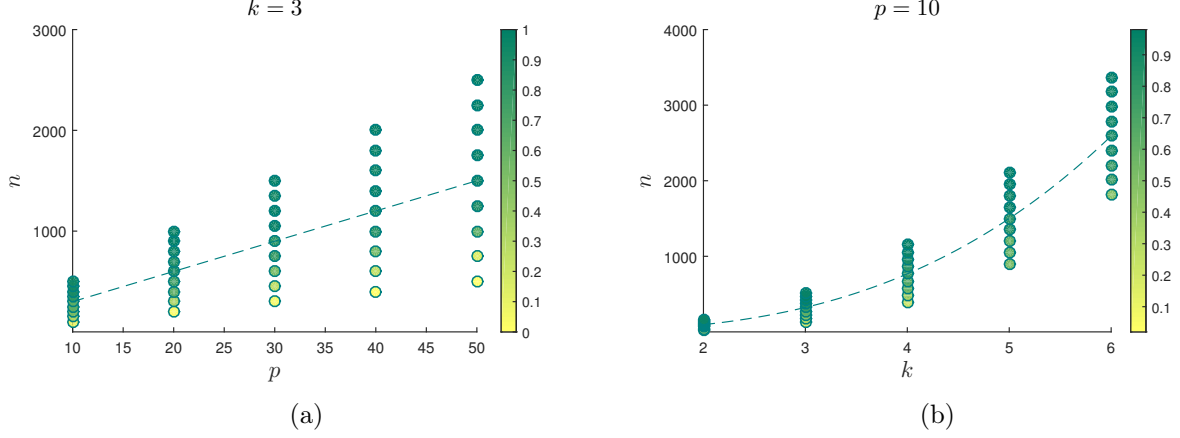
**Datasets.** For given problem size  $(n, p, k)$ , we generate synthetic datasets as follows. Covariate vectors  $\{\mathbf{x}_i\}_{i=1}^n$  are drawn independently from  $\mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ . Model parameters  $\{\beta_j\}_{j=1}^k$  are a random set of  $k$  vectors in  $\mathbb{S}^{p-1}$ , where every two distinct  $\beta$ s have distance  $\Delta = 1.2$ . Therefore, these parameters are not orthogonal. Suppose  $\mathbf{B} \in \mathbb{R}^{p \times k}$  denotes the matrix with  $\beta_j$  as the  $j$ -th column. We let  $\mathbf{B} = \mathbf{U}\mathbf{\Lambda}^{1/2}\mathbf{V}^\top$ , where  $\mathbf{U} \in \mathbb{R}^{p \times k}$  represents the basis of a random  $k$ -dimensional subspace in  $\mathbb{R}^p$ . Matrices  $\mathbf{\Lambda}, \mathbf{V} \in \mathbb{R}^{k \times k}$  are from the eigen-decomposition of symmetric matrix  $\mathbf{C} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top$ , where the diagonal terms of  $\mathbf{C}$  are 1 and the rest entries are  $1 - \Delta^2/2$ . We assign equal weights  $\omega_j = 1/k$  for all clusters.

**Results.** Our first set of results, presented in Figure 1, show the convergence of estimation errors of AltMin with random and tensor initializations. Recall that estimation error is defined in (12). In random setting, AltMin starts with a set of uniformly random  $k$  vectors in  $\mathbb{S}^{p-1}$ . We find that AltMin with random starting points has quite slow convergence, and fails to produce true  $\beta$ s with significant probability. In contrast, with the same amount of samples, tensor method provides more accurate starting points, which leads to much faster convergence of AltMin to the global optima. These results thus back up our convergence theory of AltMin (Theorem 2), and demonstrate the power of using tensor decomposition initialization.

The second set of results, presented in Figure 2, explore the statistical efficiency of the proposed algorithm—tensor initialized AltMin. For fixed  $k = 3$ , Figure (2a) reveals a linear dependence of the necessary sample size on  $p$ , which matches our results in Corollary 1. With fixed  $p$ , Figure (2b) indicates that  $\mathcal{O}(k^3)$  samples could be enough in practice, which is much better than our theoretical guarantee  $\mathcal{O}(k^{10})$ . Sharpening the polynomial factor on  $k$  is an interesting direction of future research.

## 5 Proofs

In this section, we provide proofs for Lemma 1 and the results presented in Section 3.



**Figure 2.** Exact recovery probability of “Tensor + AltMin” with varied  $(n, p, k)$ . The color of every dot represents the recovery probability computed from 100 independent trials according to the colorbar on the right side. Tensor decomposition is implemented with  $L = 200k^2, N = 20 \log(k)$ . The dashed line in (a) shows function  $n = 30p$ . The dashed line in (b) shows function  $n = 12k^3$ .

## 5.1 Proof of Lemma 1

Recall that  $z_i$  denotes the latent label associated with each sample. Suppose  $X \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ , and  $Z$  has the distribution of each  $z_i$ . We find that

$$\mathbb{E}[m_0] = \sum_{j \in [k]} \mathbb{E}[\langle X, \beta_j \rangle^2] \cdot \mathbb{P}(Z = j) = \sum_{j \in [k]} \omega_j \|\beta_j\|_2^2, \quad (15)$$

$$\mathbb{E}[m_1] = \frac{1}{6} \sum_{j \in [k]} \mathbb{E}[\langle X, \beta_j \rangle^3 X] \cdot \mathbb{P}(Z = j).$$

One can check that for any  $\beta$ ,  $\mathbb{E}[\langle X, \beta \rangle^3 X] = 3 \|\beta\|_2^2 \beta$ . Therefore,

$$\mathbb{E}[m_1] = \frac{1}{2} \sum_{j \in [k]} \omega_j \|\beta_j\|_2^2 \beta_j. \quad (16)$$

For  $M_2$ , plugging (15) into (7) yields

$$\mathbb{E}[M_2] = \frac{1}{2} \sum_{j \in [k]} \omega_j \mathbb{E}[\langle X, \beta_j \rangle^2 X \otimes X] - \frac{1}{2} \sum_{j \in [k]} \omega_j \|\beta_j\|_2^2 \cdot \mathbf{I}_p.$$

One can check  $\mathbb{E}[\langle X, \beta_j \rangle^2 X \otimes X] = 2\beta_j \beta_j^\top + \|\beta_j\|_2^2 \mathbf{I}_p$ , which leads to  $\mathbb{E}[M_2] = \sum_{j \in [k]} \omega_j \beta_j \beta_j^\top$ .

For  $M_3$ , plugging (16) into (8) gives

$$\mathbb{E}[M_3] = \frac{1}{6} \sum_{j \in [k]} \omega_j \mathbb{E}[\langle X, \beta_j \rangle^3 X^{\otimes 3}] - \frac{1}{2} \sum_{j \in [k]} \omega_j \mathcal{T}(\|\beta_j\|_2^2 \beta_j).$$

Then it remains to show that for any  $\beta$ ,

$$\mathbb{E}[\langle X, \beta \rangle^3 X^{\otimes 3}] = 6\beta^{\otimes 3} + 3\mathcal{T}(\|\beta\|_2^2 \beta). \quad (17)$$



We directly verify the above inequality. Let  $X = (X_1, \dots, X_p)^\top$ ,  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$ . For  $(i, j, k) \in [p] \times [p] \times [p]$ , let  $L_{ijk}, R_{ijk}$  be the  $(i, j, k)$ -th entries of  $\mathbb{E}[\langle X, \boldsymbol{\beta} \rangle^3 X^{\otimes 3}]$  and  $6\boldsymbol{\beta}^{\otimes 3} + 3\mathcal{T}(\|\boldsymbol{\beta}\|_2^2 \boldsymbol{\beta})$  respectively. Due to symmetry, it is sufficient to consider the following cases.

- $i \neq j \neq k \neq i$ . We have  $R_{ijk} = 6\beta_i\beta_j\beta_k$ . Meanwhile,

$$L_{ijk} = \mathbb{E}[\langle X, \boldsymbol{\beta} \rangle^3 X_i X_j X_k] = \mathbb{E}[6\beta_i\beta_j\beta_k X_i^2 X_j^2 X_k^2] = 6\beta_i\beta_j\beta_k.$$

- $i = j \neq k$ . We have  $R_{ijk} = 6\beta_i^2\beta_k + 3\|\boldsymbol{\beta}\|_2^2\beta_k$ , and

$$\begin{aligned} L_{ijk} &= \mathbb{E}[\langle X, \boldsymbol{\beta} \rangle^3 X_i^2 X_k] \\ &= \mathbb{E}[\beta_k^3 X_i^2 X_k^4] + \mathbb{E}[3\beta_i^2\beta_k X_i^4 X_k^2] + \sum_{t \in [p], t \neq i, k} \mathbb{E}[3\beta_k\beta_t^2 X_i^2 X_k^2 X_t^2] \\ &= 3\beta_k^3 + 9\beta_i^2\beta_k + 3\beta_k(\|\boldsymbol{\beta}\|_2^2 - \beta_i^2 - \beta_k^2) = 6\beta_i^2\beta_k + 3\|\boldsymbol{\beta}\|_2^2\beta_k. \end{aligned}$$

- $i = j = k$ . We have  $R_{ijk} = 6\beta_i^3 + 9\|\boldsymbol{\beta}\|_2^2\beta_i$ , and

$$\begin{aligned} L_{ijk} &= \mathbb{E}[\langle X, \boldsymbol{\beta} \rangle^3 X_i^3] = \mathbb{E}[\beta_i^3 X_i^6] + \sum_{j \in [p], j \neq i} \mathbb{E}[3\beta_i\beta_j^2 X_i^4 X_j^2] \\ &= 15\beta_i^3 + 9\beta_i(\|\boldsymbol{\beta}\|_2^2 - \beta_i^2) = 6\beta_i^2\beta_k + 3\|\boldsymbol{\beta}\|_2^2\beta_k. \end{aligned}$$

In the above calculation, we frequently used the fact that odd-order moments of symmetric Gaussian is 0. We finish proving (17), and thus conclude the proof.

## 5.2 Proof of Lemma 2

Recall that  $\sigma_k = \sigma_k(\overline{\mathbf{M}}_2) = \sigma_k(\sum_{j \in [k]} \omega_k \boldsymbol{\beta}_j \boldsymbol{\beta}_j^\top)$ . We always have

$$\sigma_k \geq \underline{\omega} \sigma_k \left( \sum_{j \in [k]} \boldsymbol{\beta}_j \boldsymbol{\beta}_j^\top \right).$$

Let  $\mathbf{B} \in \mathbb{R}^{p \times k}$  be the matrix with columns  $\boldsymbol{\beta}_j$ . Then we have

$$\sigma_k \left( \sum_{j \in [k]} \boldsymbol{\beta}_j \boldsymbol{\beta}_j^\top \right) = \sigma_{\min}(\mathbf{B}^\top \mathbf{B}).$$

Thanks to the nearly orthonormal condition, matrix  $\mathbf{D} = \mathbf{B}^\top \mathbf{B}$  has diagonal terms greater than  $1 - \eta$  and the rest entries have magnitude smaller than  $\gamma$ . Therefore, for any  $\mathbf{u} \in \mathbb{S}^{k-1}$ , we have

$$\mathbf{u}^\top \mathbf{D} \mathbf{u} \geq (1 - \eta) \|\mathbf{u}\|_2^2 - \gamma \|\mathbf{u}\|_1^2 \geq 1 - \eta - k\gamma,$$

which completes the proof.

### 5.3 Proof of Corollary 1

Linear convergence of AltMin requires  $\varepsilon_0 \lesssim (1/k^2 \wedge \underline{\omega})\Delta$ . Plugging it as accuracy into Theorem 1 shows that it suffices to let

$$n_{\text{init}} \gtrsim \frac{(k^4 + 1/\underline{\omega}^2)(p/\sigma_k^2 + k^2 + p) \log(k/\delta)}{\underline{\omega}\sigma_k^3\Delta^2} \log^3(n_{\text{init}}) + \frac{k}{\underline{\omega}\delta}.$$

The condition of  $n$  in Lemma 3 is implied by the condition of  $n/T$  in (14). Therefore, if  $\{\beta_j^{(T-1)}\}_{j=1}^k$  produced by AltMin satisfies  $\mathcal{E}(\{\beta_j^{(T-1)}\}) \leq \frac{\delta T}{4kn_{\text{alt}}}\Delta$ , Lemma 3 implies that the  $T$ -th step of AltMin (using  $n_{\text{alt}}/T$  samples) produces  $\{\beta_j\}$  with high probability. Thanks to linear convergence, we have  $\mathcal{E}(\{\beta_j^{(T-1)}\}) \leq (1/2)^{T-1}\Delta$ . So it suffices to have

$$(1/2)^{T-1}\Delta \lesssim \frac{\delta T}{kn_{\text{alt}}}\Delta.$$

Hence, choosing  $T = C \log(kn_{\text{alt}}/\delta)$  with sufficiently large constant  $C$  satisfies the above inequality. Plugging this choice of  $T$  into (14) shows that it suffices to let

$$n_{\text{alt}} \gtrsim \left( \frac{kp}{\underline{\omega}} + \frac{\log(k/\delta) + \log \log(kn_{\text{alt}}/\delta)}{\underline{\omega}^2} \right) \log(kn_{\text{alt}}/\delta).$$

Condition on  $p$  in Theorem 2 then becomes  $p \gtrsim \log(k/\delta) + \log \log(kn_{\text{alt}}/\delta)$ , under which the above requirement of  $n_{\text{alt}}$  can be strengthened to

$$n_{\text{alt}} \gtrsim \left( \frac{kp}{\underline{\omega}} + \frac{p}{\underline{\omega}^2} \right) \log(kn_{\text{alt}}/\delta).$$

### 5.4 Proofs about Tensor Decomposition

In this section, we prove the guarantee of tensor decomposition. Let  $\overline{\mathbf{M}}_2 := \mathbb{E}[\mathbf{M}_2]$  and  $\overline{\mathbf{M}}_3 := \mathbb{E}[\mathbf{M}_3]$ . The proof idea of Theorem 1 is to show how approximate the empirical moments are to their expectations, and then establish the dependence between errors of approximation and estimation. Therefore, our proofs break down into the next two subsections. In Section 5.4.1, given the approximation errors of moments

$$\varepsilon_2 := \|\mathbf{M}_2 - \overline{\mathbf{M}}_2\|_{op}, \quad (18)$$

$$\varepsilon_3 := \|\mathbf{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W}) - \overline{\mathbf{M}}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})\|_{op}, \quad (19)$$

we follow the processes shown in Algorithm 1 to obtain an upper bound of the estimation error  $\mathcal{E}(\{\beta_j^{(0)}\})$  in terms of  $\varepsilon_2$  and  $\varepsilon_3$ . In Section 5.4.2, the dependence between  $\varepsilon_2, \varepsilon_3$  and sample size is revealed by concentration analysis. We put these two parts together in Section 5.4.3 to prove Theorem 1.

#### 5.4.1 Error Transfer

We now turn to show the how error is transferred from approximation bound to initial estimation. Recall that the robust tensor power method is run on tensor  $\mathbf{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})$ . We let  $\overline{\mathbf{W}}$  be the

whitening matrix of  $\overline{\mathbf{M}}_2$ . Then tensor  $\overline{\mathbf{M}}_3(\overline{\mathbf{W}}, \overline{\mathbf{W}}, \overline{\mathbf{W}})$  has orthogonal factorization

$$\overline{\mathbf{M}}_3(\overline{\mathbf{W}}, \overline{\mathbf{W}}, \overline{\mathbf{W}}) = \sum_{j=1}^k \omega'_j \boldsymbol{\beta}'_j \otimes \boldsymbol{\beta}'_j \otimes \boldsymbol{\beta}'_j,$$

where  $\omega'_j = 1/\sqrt{\omega_j}$ ,  $\boldsymbol{\beta}'_j = \sqrt{\omega_j} \overline{\mathbf{W}}^\top \boldsymbol{\beta}_j$  and  $\|\boldsymbol{\beta}'_j\|_2 = 1$  for all  $j \in [k]$ . We will use the next theory of robust tensor power method presented in [2].

**Lemma 4** (Guarantee of Robust Tensor Power Method, Theorem 5.1 in [2]). *Suppose  $\mathbf{T} \in \mathbb{R}^{k \times k \times k}$  is a tensor with decomposition  $\mathbf{T} = \sum_{j=1}^k \lambda_j \boldsymbol{\beta}_j^{\otimes 3}$  where every  $\lambda_j > 0$  and  $\{\boldsymbol{\beta}_j\}$  are orthonormal. Put  $\bar{\lambda} := \max_{j \in [k]} \{\lambda_j\}$ ,  $\underline{\lambda} := \min_{j \in [k]} \{\lambda_j\}$ . Let  $\widehat{\mathbf{T}} = \mathbf{T} + \mathbf{E}$  be the input of Algorithm 2, where  $\mathbf{E}$  is a symmetric tensor with  $\|\mathbf{E}\|_{op} \leq \epsilon$ . There exist constants  $C_i$  such that the following holds. Suppose  $\epsilon \leq C_1 \underline{\lambda}/k$ . For any  $\delta \in (0, 1)$ , suppose  $(N, L)$  in Algorithm 2 satisfy*

$$N \geq C_2 \cdot (\log k + \log \log(\bar{\lambda}/\epsilon)), \quad L \geq C_3 \cdot \text{poly}(k) \log(1/\delta),$$

for some polynomial function  $\text{poly}(\cdot)$ . With probability at least  $1 - \delta$ ,  $\{(\widehat{\lambda}_j, \widehat{\boldsymbol{\beta}}_j)\}$  returned by Algorithm 2 satisfy the bound

$$\|\widehat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}_{\pi(j)}\|_2 \leq 8\epsilon/\lambda_{\pi(j)}, \quad |\widehat{\lambda}_j - \lambda_{\pi(j)}| \leq 5\epsilon, \quad \text{for all } j \in [k],$$

where  $\pi(\cdot)$  is some permutation function on  $[k]$ .

Without loss of generality, we set the permutation  $\pi(\cdot)$  in the above result to be identity. Lemma 4 implies that if

$$\epsilon := \|\overline{\mathbf{M}}_3(\overline{\mathbf{W}}, \overline{\mathbf{W}}, \overline{\mathbf{W}}) - \mathbf{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})\|_{op} \lesssim \frac{1}{k},$$

then with high probability,  $\{(\tilde{\omega}_j, \tilde{\boldsymbol{\beta}}_j)\}_{j=1}^k$  produced in the line 6 of Algorithm 1 satisfy

$$\|\tilde{\boldsymbol{\beta}}_j - \boldsymbol{\beta}'_j\|_2 \leq 8\epsilon/\omega'_j = 8\epsilon\sqrt{\omega_j}, \quad |\tilde{\omega}_j - \omega'_j| \leq 5\epsilon.$$

Then we have

$$\begin{aligned} \|\boldsymbol{\beta}_j^{(0)} - \boldsymbol{\beta}_j\|_2 &= \|\tilde{\omega}_j(\mathbf{W}^\top)^\dagger \tilde{\boldsymbol{\beta}}_j - \omega'_j(\overline{\mathbf{W}}^\top)^\dagger \boldsymbol{\beta}'_j\|_2 \\ &\leq \|\tilde{\omega}_j(\mathbf{W}^\top)^\dagger \tilde{\boldsymbol{\beta}}_j - \tilde{\omega}_j(\overline{\mathbf{W}}^\top)^\dagger \boldsymbol{\beta}'_j\|_2 + \|\tilde{\omega}_j(\overline{\mathbf{W}}^\top)^\dagger \boldsymbol{\beta}'_j - \omega'_j(\overline{\mathbf{W}}^\top)^\dagger \boldsymbol{\beta}'_j\|_2 \\ &\leq \|\tilde{\omega}_j(\mathbf{W}^\top)^\dagger \tilde{\boldsymbol{\beta}}_j - \tilde{\omega}_j(\overline{\mathbf{W}}^\top)^\dagger \boldsymbol{\beta}'_j\|_2 + 5\epsilon \|\overline{\mathbf{W}}^\dagger\|_{op} \\ &\leq \|\tilde{\omega}_j(\mathbf{W}^\top)^\dagger \tilde{\boldsymbol{\beta}}_j - \tilde{\omega}_j(\mathbf{W}^\top)^\dagger \boldsymbol{\beta}'_j\|_2 + \|\tilde{\omega}_j(\mathbf{W}^\top)^\dagger \boldsymbol{\beta}'_j - \tilde{\omega}_j(\overline{\mathbf{W}}^\top)^\dagger \boldsymbol{\beta}'_j\|_2 + 5\epsilon \|\overline{\mathbf{W}}^\dagger\|_{op} \\ &\leq \|\tilde{\omega}_j(\mathbf{W}^\top)^\dagger \tilde{\boldsymbol{\beta}}_j - \tilde{\omega}_j(\mathbf{W}^\top)^\dagger \boldsymbol{\beta}'_j\|_2 + \|\mathbf{W}^\dagger - \overline{\mathbf{W}}^\dagger\|_{op} + 5\epsilon \|\overline{\mathbf{W}}^\dagger\|_{op} \\ &\leq \|\mathbf{W}^\dagger\|_{op} \cdot \|\tilde{\boldsymbol{\beta}}_j - \boldsymbol{\beta}'_j\|_2 + \|\mathbf{W}^\dagger - \overline{\mathbf{W}}^\dagger\|_{op} + 5\epsilon \|\overline{\mathbf{W}}^\dagger\|_{op} \\ &\leq 8\epsilon\sqrt{\omega_j} \|\mathbf{W}^\dagger\|_{op} + \|\mathbf{W}^\dagger - \overline{\mathbf{W}}^\dagger\|_{op} + 5\epsilon \|\overline{\mathbf{W}}^\dagger\|_{op}. \end{aligned}$$

Recall that  $\sigma_k = \sigma_k(\overline{\mathbf{M}}_2)$ . Put  $\sigma_1 = \sigma_1(\overline{\mathbf{M}}_2)$ . Let  $\mathbf{M}'_2$  be the best rank  $k$  approximation of  $\mathbf{M}_2$ . We have

$$\|\mathbf{M}'_2 - \overline{\mathbf{M}}_2\|_{op} \leq \epsilon_2 + \sigma_{k+1}(\mathbf{M}_2) \leq 2\epsilon_2,$$

where the last step follows from Weyl's theorem. Using the properties of whitening in Lemma 9 by replacing  $\mathbf{A}, \widehat{\mathbf{A}}$  with  $\overline{\mathbf{M}}_2, \mathbf{M}'_2$ , when  $\epsilon_2/\sigma_k \leq 1/6$ , we have

$$\begin{aligned} \|\mathbf{W}^\dagger\|_{op} &\leq 2\|\overline{\mathbf{W}}^\dagger\|_{op} = 2\sqrt{\sigma_1}, \\ \|\mathbf{W}^\dagger - \overline{\mathbf{W}}^\dagger\|_{op} &\leq 4\epsilon_2\|\overline{\mathbf{W}}^\dagger\|_{op}/\sigma_k = 4\epsilon_2\sqrt{\sigma_1}/\sigma_k. \end{aligned}$$

We thus obtain

$$\|\beta_j^{(0)} - \beta_j\|_2 \leq 21\sqrt{\sigma_1}\epsilon + \frac{4\sqrt{\sigma_1}\epsilon_2}{\sigma_k}. \quad (20)$$

It remains to relate  $\epsilon$  to  $\epsilon_2$  and  $\epsilon_3$ . We apply a series of triangle inequalities as follows.

$$\begin{aligned} \epsilon &= \|\overline{\mathbf{M}}_3(\overline{\mathbf{W}}, \overline{\mathbf{W}}, \overline{\mathbf{W}}) - \mathbf{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})\|_{op} \\ &= \|\overline{\mathbf{M}}_3(\overline{\mathbf{W}}, \overline{\mathbf{W}}, \overline{\mathbf{W}}) - \overline{\mathbf{M}}_3(\overline{\mathbf{W}}, \overline{\mathbf{W}}, \mathbf{W}) + \overline{\mathbf{M}}_3(\overline{\mathbf{W}}, \overline{\mathbf{W}}, \mathbf{W}) - \mathbf{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})\|_{op} \\ &\leq \|\overline{\mathbf{M}}_3(\overline{\mathbf{W}}, \overline{\mathbf{W}}, \mathbf{W}) - \overline{\mathbf{M}}_3(\overline{\mathbf{W}}, \mathbf{W}, \mathbf{W}) + \overline{\mathbf{M}}_3(\overline{\mathbf{W}}, \mathbf{W}, \mathbf{W}) - \mathbf{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})\|_{op} \\ &\quad + \|\overline{\mathbf{M}}_3(\overline{\mathbf{W}}, \overline{\mathbf{W}}, \overline{\mathbf{W}} - \mathbf{W})\|_{op} \\ &\leq \|\overline{\mathbf{M}}_3\|_{op} \|\overline{\mathbf{W}}\|_{op}^2 \|\overline{\mathbf{W}} - \mathbf{W}\|_{op} + \|\overline{\mathbf{M}}_3(\overline{\mathbf{W}}, \overline{\mathbf{W}}, \mathbf{W}) - \overline{\mathbf{M}}_3(\overline{\mathbf{W}}, \mathbf{W}, \mathbf{W})\|_{op} \\ &\quad + \|\overline{\mathbf{M}}_3(\overline{\mathbf{W}}, \mathbf{W}, \mathbf{W}) - \mathbf{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})\|_{op} \\ &\leq \|\overline{\mathbf{M}}_3\|_{op} \|\overline{\mathbf{W}}\|_{op}^2 \|\overline{\mathbf{W}} - \mathbf{W}\|_{op} + \|\overline{\mathbf{M}}_3\|_{op} \|\overline{\mathbf{W}}\|_{op} \|\mathbf{W}\|_{op} \|\overline{\mathbf{W}} - \mathbf{W}\|_{op} \\ &\quad + \|\overline{\mathbf{M}}_3\|_{op} \|\mathbf{W}\|_{op}^2 \|\overline{\mathbf{W}} - \mathbf{W}\|_{op} + \epsilon_3. \end{aligned} \quad (21)$$

Applying Lemma 9 again, we have

$$\|\mathbf{W}\|_{op} \leq 2\|\overline{\mathbf{W}}\|_{op} = 2/\sqrt{\sigma_k}, \quad \|\overline{\mathbf{W}} - \mathbf{W}\|_{op} \leq 4\epsilon_2/\sqrt{\sigma_k^3}. \quad (22)$$

Plugging it back into the last line of (21) yields

$$\|\overline{\mathbf{M}}_3(\overline{\mathbf{W}}, \overline{\mathbf{W}}, \overline{\mathbf{W}}) - \mathbf{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})\|_{op} \leq \frac{28\epsilon_2}{\sqrt{\sigma_k^5}} \|\overline{\mathbf{M}}_3\|_{op} + \epsilon_3. \quad (23)$$

We thus obtain the following error bound by putting (20) and (23) together:

$$\|\beta_j^{(0)} - \beta_j\|_2 \lesssim \frac{\sqrt{\sigma_1}\epsilon_2}{\sigma_k} + \sqrt{\sigma_1}\epsilon_3 + \frac{\sqrt{\sigma_1}\|\overline{\mathbf{M}}_3\|_{op}\epsilon_2}{\sqrt{\sigma_k^5}}, \quad \text{for all } j \in [k]. \quad (24)$$

Recall that, in order to obtain (20), we have to make sure  $\epsilon \lesssim 1/k$  as required in Lemma 4. Then inequality (23) indicates that it's sufficient to require

$$\epsilon_2 \lesssim \frac{\sqrt{\sigma_k^5}}{k\|\overline{\mathbf{M}}_3\|_{op}} \quad \text{and} \quad \epsilon_3 \lesssim \frac{1}{k}, \quad (25)$$

which will be used in the concentration analysis.

### 5.4.2 Concentration Analysis

Now we turn to the analysis of the concentration of empirical moments, and we derive upper bounds on  $\epsilon_2$  and  $\epsilon_3$ . Note that  $\mathbf{M}_3$  involves Gaussian's high-order moments (up to 6th moment). In order to deal with the heavy tail, we will leverage a *truncation argument*, where we introduce truncated response  $y'_i$  as

$$y'_i = \begin{cases} y_i, & \text{if } |y_i| \leq M, \\ \text{sign}(y_i) \cdot M, & \text{otherwise} \end{cases}, \quad (26)$$

where  $M > 0$  is some threshold chosen in our analysis. When  $M$  is sufficiently large, we have  $y_i = y'_i$  for all  $i \in [n]$  with high probability, which means the tail bounds about  $\{(y'_i, \mathbf{x}_i)\}$  still apply to original samples  $\{(y_i, \mathbf{x}_i)\}$ . The advance of analyzing concentration using  $(y'_i, \mathbf{x}_i)$  is that  $y'_i \cdot \mathbf{x}_i$  is sub-Gaussian random vector thanks to the boundedness of  $y'_i$ . One should note that truncating  $y_i$  might change the expectation of moments slightly. Therefore, a tedious but important part of our analysis is to show that the expectation deviation from truncation is much smaller compared to the desired tail bound. In detail, we have the next result proved using the truncation idea. See Section 6.1 for the complete proof.

**Lemma 5** (Concentration of Empirical Moments of Single Model). *Suppose  $n$  samples  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are generated from  $\mathcal{N}(\mathbf{0}, \mathbf{I}_p)$  and  $y_i = \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle$  for some fixed  $\boldsymbol{\beta} \in \mathbb{S}^{p-1}$ . Let*

$$\mathbf{m}_1 = \frac{1}{n} \sum_{i \in [n]} y_i^3 \mathbf{x}_i, \quad \mathbf{M}_2 = \frac{1}{n} \sum_{i \in [n]} y_i^2 \mathbf{x}_i^{\otimes 2}, \quad \mathbf{M}_3 = \frac{1}{n} \sum_{i \in [n]} y_i^3 \mathbf{x}_i^{\otimes 3}.$$

Moreover, let  $\bar{m}_0 = \mathbb{E}[m_0], \bar{\mathbf{m}}_1 = \mathbb{E}[\mathbf{m}_1], \bar{\mathbf{M}}_2 = \mathbb{E}[\mathbf{M}_2], \bar{\mathbf{M}}_3 = \mathbb{E}[\mathbf{M}_3]$ . There exist constants  $C_i$  such that the following holds. Pick any  $\delta \in (0, 1)$  and any fixed matrix  $\mathbf{S} \in \mathbb{R}^{p \times s}$  with  $s \leq p$ .

1. If  $n \geq C_1/\delta$ , with probability at least  $1 - \delta$ , we have

$$\left\| \mathbf{S}^\top (\mathbf{m}_1 - \bar{\mathbf{m}}_1) \right\|_2 \leq C_2 \|\mathbf{S}\|_{op} \frac{\log^{3/2} n}{\sqrt{n}} \max \left\{ \sqrt{\log \left( \frac{2}{\delta} \right)}, \sqrt{s} \right\}. \quad (27)$$

2. If  $n \geq C_3 \max\{1/\delta, s\}$ , with probability  $1 - \delta$ , we have

$$\left\| \mathbf{S}^\top (\mathbf{M}_2 - \bar{\mathbf{M}}_2) \mathbf{S} \right\|_{op} \leq C_4 \|\mathbf{S}\|_{op}^2 \frac{\log n}{\sqrt{n}} \max \left\{ \sqrt{\log \left( \frac{2}{\delta} \right)}, \sqrt{s} \right\}. \quad (28)$$

3. If  $n \geq C_5 \max\{s \log \left( \frac{2}{\delta} \right), 1/\delta\}$ , with probability at least  $1 - \delta$ , we have

$$\left\| (\mathbf{M}_3 - \bar{\mathbf{M}}_3) (\mathbf{S}, \mathbf{S}, \mathbf{S}) \right\|_{op} \leq C_6 \|\mathbf{S}\|_{op}^3 \frac{s \log^{3/2} n}{\sqrt{n}} \sqrt{\log \left( \frac{2}{\delta} \right)}. \quad (29)$$

This result provides concentration bounds of the moments constructed from single linear model. In the case of mixture samples  $\{(y_i, \mathbf{x}_i)\}_{i=1}^n$ , we can split the set into  $k$  sets  $\{(y_i^{(j)}, \mathbf{x}_i^{(j)})\}_{i=1}^{n_j}, j =$

$1, \dots, k$ , where the  $j$ -th set corresponds to linear model  $\beta_j$ . Therefore, for the moments given in (4)-(6), we have

$$\begin{aligned} m_0 &= \sum_{j \in [k]} \hat{\omega}_j m_0^{(j)}, \quad \mathbf{m}_1 = \frac{1}{6} \sum_{j \in [k]} \hat{\omega}_j \mathbf{m}_1^{(j)}, \\ \mathbf{M}_2 &= \frac{1}{2} \sum_{j \in [k]} \hat{\omega}_j \mathbf{M}_2^{(j)} - \frac{1}{2} m_0 \cdot \mathbf{I}_p, \quad \mathbf{M}_3 = \frac{1}{6} \sum_{j \in [k]} \hat{\omega}_j \mathbf{M}_3^{(j)} - \mathcal{T}(\mathbf{m}_1), \end{aligned}$$

where  $\hat{\omega}_j$  denotes the empirical proportion of each model, and we let  $m_0^{(j)} := \frac{1}{n_j} \sum_{i \in [n_j]} y_i^{(j)2}$ ,  $\mathbf{m}_1^{(j)} := \frac{1}{n_j} \sum_{i \in [n_j]} y_i^{(j)2} \mathbf{x}_i^{(j)}$ ,  $\mathbf{M}_2^{(j)} := \frac{1}{n_j} \sum_{i \in [n_j]} y_i^{(j)2} \mathbf{x}_i^{(j)\otimes 2}$ ,  $\mathbf{M}_3^{(j)} := \frac{1}{n_j} \sum_{i \in [n_j]} y_i^{(j)3} \mathbf{x}_i^{(j)\otimes 3}$ .

Next, we will derive concentration bounds for  $m_0, \mathbf{m}_1, \mathbf{M}_2, \mathbf{M}_3$  respectively. To ease notation, for every moment, we use  $n$  to denote the number of samples for computing it, while they might be computed from different sets of samples in Algorithm 1.

**Bound of  $|m_0 - \bar{m}_0|$ .** We find that

$$\begin{aligned} \epsilon_0 &:= |m_0 - \bar{m}_0| \leq \sum_{j \in [k]} \hat{\omega}_j \left| m_0^{(j)} - \mathbb{E}[m_0^{(j)}] \right| + \sum_{j \in [k]} |\hat{\omega}_j - \omega_j| \cdot \mathbb{E}[m_0^{(j)}] \\ &\leq \sup_{j \in [k]} |m_0^{(j)} - \mathbb{E}[m_0^{(j)}]| + \sum_{j \in [k]} |\hat{\omega}_j - \omega_j| \cdot \mathbb{E}[m_0^{(j)}] \\ &\leq \sup_{j \in [k]} |m_0^{(j)} - \mathbb{E}[m_0^{(j)}]| + \underbrace{\sum_{j \in [k]} |\hat{\omega}_j - \omega_j|}_{\epsilon_\omega}, \end{aligned} \tag{30}$$

where the last step follows from the fact  $\mathbb{E}[m_0^{(j)}] \leq 1$  due to the assumption  $\max_{j \in [k]} \|\beta_j\|_2 = 1$ . We first bound  $\epsilon_\omega$ . Note that  $n\hat{\omega}_j$  is a sum of  $n$  Bernoulli random variables with success probability  $\omega_j$ . Lemma 8 gives that for any  $t \in (0, 1)$

$$\mathbb{P}(|\hat{\omega}_j - \omega_j| \geq t\omega_j) \leq 2e^{-\frac{3t^2}{2(t+3)}n\omega_j} \leq 2e^{-3t^2n\omega/8}.$$

Using union bound and setting  $t = \sqrt{8 \log(2k/\delta)/(3\omega n)}$ , which can be less than 1 when  $n \geq C \log(k/\delta)/\omega$  for sufficiently large  $C$ , we obtain

$$\mathbb{P}\left(\epsilon_\omega \geq \sqrt{\frac{8 \log(2k/\delta)}{3\omega n}}\right) \leq 2ke^{-3t^2n\omega/8} = \delta. \tag{31}$$

Now we turn to the first term in (30). Note that  $y_i^2 = \langle \mathbf{x}_i, \beta_j \rangle^2$  is sub-Gaussian with constant Orlicz norm as  $\|\beta_j\|_2 \leq 1$ . Then by standard concentration of sub-Gaussian (e.g., (59) with  $p = 1$ ), we find that there exist constants  $C, C'$  such that if  $n \geq C \frac{1}{\omega} \log(k/\delta)$ , we have

$$\mathbb{P}\left(\sup_{j \in [k]} |m_0^{(j)} - \mathbb{E}[m_0^{(j)}]| \geq C' \sqrt{\frac{1}{\omega n} \log\left(\frac{k}{\delta}\right)}\right) \leq \delta.$$

for any  $\delta \in (0, 1)$ . Excluding the probability  $\delta$ , we obtain

$$\epsilon_0 \lesssim \sqrt{\log(k/\delta)/(\omega n)} + \epsilon_\omega. \tag{32}$$

**Bound of  $\|\mathbf{m}_1 - \overline{\mathbf{m}}_1\|_2$ .** Similar to (30), we have

$$\begin{aligned}\epsilon_1 &:= \|\mathbf{m}_1 - \overline{\mathbf{m}}_1\|_2 \lesssim \sup_{j \in [k]} \left\| \mathbf{m}_1^{(j)} - \mathbb{E}[\mathbf{m}_1^{(j)}] \right\|_2 + \epsilon_\omega \cdot \sup_{j \in [k]} \left\| \mathbb{E}[\mathbf{m}_1^{(j)}] \right\|_2 \\ &\lesssim \sup_{j \in [k]} \left\| \mathbf{m}_1^{(j)} - \mathbb{E}[\mathbf{m}_1^{(j)}] \right\|_2 + \epsilon_\omega.\end{aligned}$$

Using (27) in Lemma 5 by setting  $\mathbf{S} = \mathbf{I}_p$  and replacing  $\delta$  with  $\delta/k$ , we have that the condition  $n \gtrsim k/(\underline{\omega}\delta)$  leads to

$$\sup_{j \in [k]} \left\| \mathbf{m}_1^{(j)} - \mathbb{E}[\mathbf{m}_1^{(j)}] \right\|_2 \lesssim \frac{\log^{3/2}(\underline{\omega}n)}{\sqrt{\underline{\omega}n}} \sqrt{p \log(2k/\delta)}$$

holds with probability at least  $1 - \delta$ . Conditioning on this event leads to

$$\epsilon_1 \lesssim \log^{3/2}(\underline{\omega}n) \sqrt{p \log(2k/\delta)/(\underline{\omega}n)} + \epsilon_\omega. \quad (33)$$

**Bound of  $\|\mathbf{M}_2 - \overline{\mathbf{M}}_2\|_{op}$ .** We find that

$$\begin{aligned}\epsilon_2 &\lesssim \left\| \sum_{j \in [k]} \widehat{\omega}_j \mathbf{M}_2^{(j)} - \sum_{j \in [k]} \omega_j \mathbb{E}[\mathbf{M}_2^{(j)}] \right\|_{op} + |m_0 - \overline{m}_0| \\ &\lesssim \sup_{j \in [k]} \left\| \mathbf{M}_2^{(j)} - \mathbb{E}[\mathbf{M}_2^{(j)}] \right\|_{op} + \epsilon_\omega + \epsilon_0,\end{aligned}$$

where the second step follows from similar calculation in (30) and the fact  $\left\| \mathbb{E}[\mathbf{M}_2^{(j)}] \right\|_{op} \lesssim 1$  for all  $j \in [k]$ . Applying (28) by choosing  $\mathbf{S} = \mathbf{I}_p$  and setting  $\delta$  to be  $\delta/k$ , we have that when  $n \gtrsim \underline{\omega}^{-1} \max\{k/\delta, p\}$ ,

$$\sup_{j \in [k]} \left\| \mathbf{M}_2^{(j)} - \mathbb{E}[\mathbf{M}_2^{(j)}] \right\|_{op} \lesssim \log(\underline{\omega}n) \sqrt{p \log(2k/\delta)/(\underline{\omega}n)}$$

holds with probability at least  $1 - \delta$ . Conditioning on the event, we conclude that

$$\epsilon_2 \lesssim \log(\underline{\omega}n) \sqrt{p \log(2k/\delta)/(\underline{\omega}n)} + \epsilon_\omega + \epsilon_0. \quad (34)$$

**Bound of  $\|\mathbf{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W}) - \overline{\mathbf{M}}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})\|_{op}$ .** Now we condition on the event  $\epsilon_2/\sigma_k < 1/6$ , which can lead to  $\|\mathbf{W}\|_{op} \leq 2/\sqrt{\sigma_k}$  as shown in (22). Let  $\epsilon_{\mathcal{T}} := \|\mathcal{T}(\mathbf{m}_1 - \overline{\mathbf{m}}_1)(\mathbf{W}, \mathbf{W}, \mathbf{W})\|_{op}$ . Recall that  $\epsilon_3$  is defined in (19). We find that

$$\begin{aligned}\epsilon_3 &\lesssim \left\| \sum_{j \in [k]} \widehat{\omega}_j \mathbf{M}_3^{(j)}(\mathbf{W}, \mathbf{W}, \mathbf{W}) - \sum_{j \in [k]} \omega_j \mathbb{E}[\mathbf{M}_3^{(j)}(\mathbf{W}, \mathbf{W}, \mathbf{W})] \right\|_{op} + \epsilon_{\mathcal{T}} \\ &\lesssim \sup_{j \in [k]} \left\| \mathbf{M}_3^{(j)}(\mathbf{W}, \mathbf{W}, \mathbf{W}) - \mathbb{E}[\mathbf{M}_3^{(j)}(\mathbf{W}, \mathbf{W}, \mathbf{W})] \right\|_{op} + \epsilon_\omega + \epsilon_{\mathcal{T}}.\end{aligned}$$

Again, the last step follows from similar steps in (30) and the fact that

$$\left\| \mathbb{E}[\mathbf{M}_3^{(j)}(\mathbf{W}, \mathbf{W}, \mathbf{W})] \right\|_{op} \lesssim \left\| \mathbb{E}[\mathbf{M}_3^{(j)}] \right\|_{op} \cdot \|\mathbf{W}\|_{op}^3 \lesssim \|\mathbf{W}\|_{op}^3 \lesssim 1/\sqrt{\sigma_k^3}.$$

Note that  $\mathbf{W}$  is computed from  $\mathbf{M}_2$ . Due to the sample splitting in Algorithm 1,  $\mathbf{W}$  is independent of  $\mathbf{M}_3$ . Therefore, we can apply (29) by replacing  $\mathbf{S}, \delta$  with  $\mathbf{W}, \delta/k$  to obtain that

$$\sup_{j \in [k]} \left\| \mathbf{M}_3^{(j)}(\mathbf{W}, \mathbf{W}, \mathbf{W}) - \mathbb{E}[\mathbf{M}_3^{(j)}(\mathbf{W}, \mathbf{W}, \mathbf{W})] \right\|_{op} \lesssim \frac{1}{\sqrt{\sigma_k^3}} k \log^{3/2}(\underline{\omega}n) \sqrt{\log(2k/\delta)/(\underline{\omega}n)} \quad (35)$$

holds with probability at least  $1 - \delta$  under condition  $n \gtrsim k/(\underline{\omega}\delta)$ . For  $\mathcal{T}(\cdot)$ , we have that for any  $\mathbf{u} \in \mathbb{R}^p$ ,

$$\|\mathcal{T}(\mathbf{u})\|_{op} \leq 3 \|\mathbf{u}\|_2, \quad (36)$$

which is proved at the end of this section. We have

$$\epsilon_{\mathcal{T}} \lesssim \|\mathbf{W}\|_{op}^3 \epsilon_1 \lesssim \epsilon_1/\sqrt{\sigma_k^3}.$$

Conditioning on (35) leads to

$$\epsilon_3 \lesssim \frac{1}{\sqrt{\sigma_k^3}} k \log^{3/2}(\underline{\omega}n) \sqrt{\log(2k/\delta)/(\underline{\omega}n)} + \epsilon_{\omega} + \epsilon_1/\sqrt{\sigma_k^3}. \quad (37)$$

*Proof of Inequality (36).* For any  $\mathbf{v} \in \mathbb{S}^{p-1}$ , we have

$$\mathcal{T}(\mathbf{u})(\mathbf{v}, \mathbf{v}, \mathbf{v}) = 3 \sum_{i,j \in [p]} u_i v_i v_j^2 = 3 \sum_{i \in [p]} u_i v_i \|\mathbf{v}\|_2^2 = 3 \langle \mathbf{u}, \mathbf{v} \rangle \leq 3 \|\mathbf{u}\|_2.$$

□

### 5.4.3 Proof of Theorem 1

With the previous analysis, we are ready to prove Theorem 1. In the first place, we combine the ingredients in Section 5.4.2. Recall that we split  $n$  samples into two parts with size  $n_1$  and  $n_2$  for computing  $m_0, \mathbf{M}_2$  and  $\mathbf{m}_1, \mathbf{M}_3$  respectively. Putting (31), (32), (34) together and using union bound, we have

$$\mathbb{P} \left( \epsilon_2 \lesssim \log(\underline{\omega}n_1) \sqrt{\frac{p \log(12k/\delta)}{\underline{\omega}n_1}} \right) \geq \delta/2 \quad (38)$$

under condition  $n_1 \gtrsim \frac{1}{\underline{\omega}} (\frac{k}{\delta} \vee p)$ . Putting (31), (33) and (37) together leads to

$$\mathbb{P} \left( \epsilon_3 \lesssim \frac{(k \vee \sqrt{p}) \log^{3/2}(\underline{\omega}n_2)}{\sqrt{\sigma_k^3}} \sqrt{\frac{\log(12k/\delta)}{\underline{\omega}n_2}} \right) \geq \delta/2 \quad (39)$$

under conditions  $n_2 \gtrsim k/(\underline{\omega}\delta)$  and  $\epsilon_2 < \sigma_k/6$ . In order to guarantee  $\left\| \boldsymbol{\beta}_j^{(0)} - \boldsymbol{\beta}_j \right\|_2 \lesssim \varepsilon$  for all  $j \in [k]$ , using the error transfer inequality (24) and noting that  $\sigma_1 \leq 1, \|\overline{\mathbf{M}}_3\|_{op} \leq 1$  under assumption  $\max_{j \in [k]} \|\boldsymbol{\beta}_j\|_2 = 1$ , it is sufficient to require

$$\epsilon_2 \lesssim \sqrt{\sigma_k^5} \varepsilon, \quad \epsilon_3 \lesssim \varepsilon. \quad (40)$$



The above condition on  $\epsilon_2$  leads to  $\epsilon_2 < \sigma_k/6$  for  $\epsilon \leq 1$ . In addition, in order to let (24) hold,  $\epsilon_2, \epsilon_3$  have to satisfy condition (25). This is implied by (40) when  $\epsilon \lesssim 1/k$ . Using the relationship between  $\epsilon_2, \epsilon_3$  and  $n_1, n_2$  in (38) and (39), it is sufficient to require

$$n_1 \gtrsim \frac{p \log(12k/\delta) \log^2(n_1)}{\underline{\omega} \sigma_k^5 \epsilon^2} \vee \frac{k}{\underline{\omega} \delta} \quad \text{and} \quad n_2 \gtrsim \frac{(k^2 \vee p) \log(12k/\delta) \log^3(n_2)}{\underline{\omega} \sigma_k^3 \epsilon^2} \vee \frac{k}{\underline{\omega} \delta},$$

which concludes our proof.

## 5.5 Proof of Alternating Minimization (Theorem 2)

It is sufficient to show the linear error decay in one step. Then the error bound for each step  $t$  can be obtained by induction. Without loss of generality, we focus on the first step  $t = 0$ . Also we assume  $\pi(j) = j$  for simplicity. Let  $B = n/T$  be the sample size in the first step. Let  $\mathcal{A}_j$  denote the index set of samples that are clustered to model  $j$  in the label assignment step, namely

$$\mathcal{A}_j := \left\{ i \in [B] \mid |y_i - \langle \mathbf{x}_i, \boldsymbol{\beta}_j^{(0)} \rangle| < |y_i - \langle \mathbf{x}_i, \boldsymbol{\beta}_t^{(0)} \rangle| \text{ for all } t \neq j \right\}.$$

We use  $\mathcal{A}_j^*$  to denote the set of samples that are truly generated from model  $\boldsymbol{\beta}_j$ , namely

$$\mathcal{A}_j^* := \{ i \in [B] \mid y_i = \langle \mathbf{x}_i, \boldsymbol{\beta}_j \rangle \}.$$

Introduce  $\epsilon_0$  as a shorthand for  $\mathcal{E}(\{\boldsymbol{\beta}_j^{(0)}\})$ . According to our assumption,  $\epsilon_0 \lesssim \Delta/k^2$ .

Let  $\boldsymbol{\Sigma}_j := \sum_{i \in \mathcal{A}_j} \mathbf{x}_i \mathbf{x}_i^\top$  be the empirical covariance of samples in  $\mathcal{A}_j$ . The updated estimate  $\boldsymbol{\beta}_j^{(1)}$  has the form

$$\boldsymbol{\beta}_j^{(1)} = \boldsymbol{\Sigma}_j^{-1} \left( \sum_{i \in \mathcal{A}_j} y_i \mathbf{x}_i \right).$$

We thus obtain

$$\begin{aligned} \boldsymbol{\beta}_j^{(1)} - \boldsymbol{\beta}_j &= \boldsymbol{\Sigma}_j^{-1} \left( \sum_{i \in \mathcal{A}_j} y_i \mathbf{x}_i \right) - \boldsymbol{\beta}_j = \boldsymbol{\Sigma}_j^{-1} \left( \sum_{i \in \mathcal{A}_j} y_i \mathbf{x}_i - \mathbf{x}_i \mathbf{x}_i^\top \boldsymbol{\beta}_j \right) \\ &= \boldsymbol{\Sigma}_j^{-1} \left( \sum_{t \in [k]} \sum_{i \in \mathcal{A}_t^* \cap \mathcal{A}_j} \mathbf{x}_i \mathbf{x}_i^\top (\boldsymbol{\beta}_t - \boldsymbol{\beta}_j) \right). \end{aligned}$$

By the Cauchy-Schwartz inequality, we obtain

$$\begin{aligned} \left\| \boldsymbol{\beta}_j^{(1)} - \boldsymbol{\beta}_j \right\|_2 &\leq \left\| \boldsymbol{\Sigma}_j^{-1} \right\|_{op} \cdot \left\| \sum_{t \in [k]} \sum_{i \in \mathcal{A}_t^* \cap \mathcal{A}_j} \mathbf{x}_i \mathbf{x}_i^\top (\boldsymbol{\beta}_t - \boldsymbol{\beta}_j) \right\|_2 \\ &\leq \underbrace{\left\| \boldsymbol{\Sigma}_j^{-1} \right\|_{op}}_{U_1} \cdot \underbrace{\left( \sum_{t \in [k]} \left\| \sum_{i \in \mathcal{A}_t^* \cap \mathcal{A}_j} \mathbf{x}_i \mathbf{x}_i^\top (\boldsymbol{\beta}_t - \boldsymbol{\beta}_j) \right\|_2 \right)}_{U_2} \end{aligned}$$

Next we bound the two terms  $U_1$  and  $U_2$  respectively.

**Bound of  $U_1$ .** First note that  $\|\boldsymbol{\Sigma}_j^{-1}\|_{op} = 1/\sigma_{\min}(\boldsymbol{\Sigma}_j)$ . We find that

$$\sigma_{\min}(\boldsymbol{\Sigma}_j) = \sigma_{\min}\left(\sum_{i \in \mathcal{A}_j} \mathbf{x}_i \mathbf{x}_i^\top\right) \geq \sigma_{\min}\left(\sum_{i \in \mathcal{A}_j \cap \mathcal{A}_j^*} \mathbf{x}_i \mathbf{x}_i^\top\right).$$

For  $X \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ , we define event  $\mathcal{E}_j$  as

$$\mathcal{E}_j := \left\{ |\langle X, \boldsymbol{\beta}_j^{(0)} - \boldsymbol{\beta}_j \rangle| \leq |\langle X, \boldsymbol{\beta}_t^{(0)} - \boldsymbol{\beta}_j \rangle|, \text{ for all } t \neq j \right\}.$$

Accordingly, we have

$$\mathbb{E}[\mathbf{x}_i \mathbf{x}_i^\top | i \in \mathcal{A}_j \cap \mathcal{A}_j^*] = \mathbb{E}[X X^\top | \mathcal{E}_j]. \quad (41)$$

To provide a lower bound of  $\mathbb{P}(\mathcal{E}_j)$ , we have

$$\begin{aligned} \mathbb{P}(\mathcal{E}_j) &= 1 - \mathbb{P}(\mathcal{E}_j^c) \geq 1 - \sum_{t \neq j} \mathbb{P}\{\langle X, \boldsymbol{\beta}_j^{(0)} - \boldsymbol{\beta}_j \rangle^2 \geq \langle X, \boldsymbol{\beta}_t^{(0)} - \boldsymbol{\beta}_j \rangle^2\} \\ &\stackrel{(a)}{\geq} 1 - \sum_{t \neq j} \frac{\|\boldsymbol{\beta}_j^{(0)} - \boldsymbol{\beta}_j\|_2}{\|\boldsymbol{\beta}_t^{(0)} - \boldsymbol{\beta}_j\|_2} \stackrel{(b)}{\geq} 1 - (k-1) \frac{4}{7k^2} \geq 1 - \frac{4}{7k}. \end{aligned} \quad (42)$$

Step (b) holds because since for all  $t \neq j$ ,

$$\frac{\|\boldsymbol{\beta}_j^{(0)} - \boldsymbol{\beta}_j\|_2}{\|\boldsymbol{\beta}_t^{(0)} - \boldsymbol{\beta}_j\|_2} \leq \frac{\|\boldsymbol{\beta}_j^{(0)} - \boldsymbol{\beta}_j\|_2}{\|\boldsymbol{\beta}_t - \boldsymbol{\beta}_j\|_2 - \|\boldsymbol{\beta}_t^{(0)} - \boldsymbol{\beta}_t\|_2} \leq \frac{\varepsilon_0}{\Delta - \varepsilon_0} \leq \frac{4}{7k^2},$$

where the last step follows from condition  $\varepsilon_0 \leq \Delta/(7k^2)$ . Step (a) in (42) is from the next result, which is proved in Section 6.2.

**Lemma 6.** Let  $X \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ . For any two fixed vectors  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^p$ , we define

$$\mathcal{E} := \{|\langle X, \mathbf{u} \rangle| \leq |\langle X, \mathbf{v} \rangle|\}.$$

We have that when  $\|\mathbf{u}\|_2 > \|\mathbf{v}\|_2$ ,

$$\mathbb{P}(\mathcal{E}) \leq \frac{\|\mathbf{v}\|_2}{\|\mathbf{u}\|_2}.$$

The next result, proved in Section 6.3, establishes the spectral structure of the covariance matrix of  $X \mid \mathcal{E}_j$ .

**Lemma 7 (Conditional Spectral Structure).** Let  $X \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ . For any  $k$  fixed vectors  $\mathbf{u}_1, \dots, \mathbf{u}_k \in \mathbb{R}^p$ , we define event

$$\mathcal{E} := \{|\langle X, \mathbf{u}_1 \rangle| \leq |\langle X, \mathbf{u}_j \rangle|, \text{ for all } j \in [k]\}.$$

When  $\mathbb{P}(\mathcal{E}) > 0$ , we have

$$\sigma_{\max}\left(\mathbb{E}[X X^\top \mid \mathcal{E}]\right) \leq k.$$

and

$$\sigma_{\min}\left(\mathbb{E}[X X^\top \mid \mathcal{E}]\right) \geq \frac{1 - k(1 - \mathbb{P}(\mathcal{E}))}{\mathbb{P}(\mathcal{E})}. \quad (43)$$

The above result suggests that

$$\sigma_{\min}(\mathbb{E}[XX^\top | \mathcal{E}_j]) \geq \frac{1 - k \cdot [1 - \mathbb{P}(\mathcal{E}_j)]}{\mathbb{P}(\mathcal{E}_j)} \geq \frac{3}{7\mathbb{P}(\mathcal{E}_j)} \geq \frac{3}{7}. \quad (44)$$

Next we will show  $\sigma_{\min}(\sum_{i \in \mathcal{A}_j \cap \mathcal{A}_j^*} \mathbf{x}_i \mathbf{x}_i^\top)$  is close to its expected value  $|\mathcal{A}_j \cap \mathcal{A}_j^*| \sigma_{\min}(\mathbb{E}[XX^\top | \mathcal{E}_j])$ . First, we prove  $|\mathcal{A}_j \cap \mathcal{A}_j^*|$  is large enough. As  $\mathbb{P}(\mathcal{E}_j) \geq 1 - 4/(7k^2) \geq 1/2$ , we have  $\mathbb{E}[|\mathcal{A}_j \cap \mathcal{A}_j^*|] \geq \mathbb{E}[\frac{1}{2}|\mathcal{A}_j^*|] = \frac{1}{2}\omega_j B$ . Therefore,  $|\mathcal{A}_j \cap \mathcal{A}_j^*|$  is summation of  $B$  independent Bernoulli random variable with success probability at least  $\omega_j/2$ . Then we have

$$\mathbb{P}\left(|\mathcal{A}_j \cap \mathcal{A}_j^*| \leq \frac{1}{4}\omega_j B\right) \leq \mathbb{P}\left(\left||\mathcal{A}_j \cap \mathcal{A}_j^*| - \mathbb{E}[|\mathcal{A}_j \cap \mathcal{A}_j^*|]\right| \geq \frac{1}{4}\omega_j B\right) \leq 2e^{-C\omega_j B} \leq 2e^{-C\omega B}, \quad (45)$$

where the second step follows from Lemma 8 and  $C$  is some constant. Conditioning on the event  $|\mathcal{A}_j \cap \mathcal{A}_j^*| \geq \omega_j B/4$ , we obtain  $|\mathcal{A}_j \cap \mathcal{A}_j^*| \gtrsim p$  when  $B \gtrsim p/\omega$ .

Note that  $X$  is sub-Gaussian random vector. Part (a) of Lemma 15 shows that  $X$  is still sub-Gaussian vector conditioning on  $\mathcal{E}_j$ . Using the conclusion  $|\mathcal{A}_j \cap \mathcal{A}_j^*| \gtrsim p$ , concentration result of sub-Gaussian in (59) (setting  $t = 1/7$  and  $K$  to be a constant) yields that, for some constant  $C$ ,

$$\mathbb{P}\left(\left\|\frac{1}{|\mathcal{A}_j \cap \mathcal{A}_j^*|} \sum_{i \in \mathcal{A}_j \cap \mathcal{A}_j^*} \mathbf{x}_i \mathbf{x}_i^\top - \mathbb{E}[XX^\top | \mathcal{E}_j]\right\|_{op} \geq \frac{1}{7}\right) \leq 2e^{-C\omega_j B} \leq 2e^{-C\omega B}. \quad (46)$$

Putting (45) and (46) together and using Weyl's theorem, we have that with probability at least  $1 - 4e^{-C'\omega B}$ ,

$$\sigma_{\min}\left(\sum_{i \in \mathcal{A}_j \cap \mathcal{A}_j^*} \mathbf{x}_i \mathbf{x}_i^\top\right) \geq |\mathcal{A}_j \cap \mathcal{A}_j^*| \cdot \left(\sigma_{\min}(\mathbb{E}[XX^\top | \mathcal{E}_j]) - \frac{1}{7}\right) \geq \frac{1}{4}\omega_j B \cdot \frac{2}{7} = \frac{1}{14}\omega_j B,$$

We thus obtain

$$\mathbb{P}(U_1 \geq 14/(w_j B)) \leq 4e^{-C'\omega B}. \quad (47)$$

**Bound of  $U_2$ .** Recall that

$$U_2 = \sum_{t \neq j} \left\| \sum_{\mathcal{A}_t^* \cap \mathcal{A}_j} \mathbf{x}_i \mathbf{x}_i^\top (\beta_t - \beta_j) \right\|_2.$$

We will bound every term with different  $t$  separately. Note that for any vector  $\mathbf{x} \in \mathbb{R}^p$  and positive semidefinite matrix  $\mathbf{A} \in \mathbb{R}^{p \times p}$ , we have

$$\|\mathbf{A}\mathbf{x}\|_2^2 \leq \sigma_{\max}(\mathbf{A})\mathbf{x}^\top \mathbf{A}\mathbf{x}.$$

Introduce  $\mathbf{Q}_t = \sum_{\mathcal{A}_i^* \cap \mathcal{A}_j} \mathbf{x}_i \mathbf{x}_i^\top$ . We find

$$\begin{aligned}
\left\| \sum_{\mathcal{A}_i^* \cap \mathcal{A}_j} \mathbf{x}_i \mathbf{x}_i^\top (\boldsymbol{\beta}_t - \boldsymbol{\beta}_j) \right\|_2^2 &\leq \sigma_{\max}(\mathbf{Q}_t) \sum_{i \in \mathcal{A}_i^* \cap \mathcal{A}_j} (\boldsymbol{\beta}_t - \boldsymbol{\beta}_j)^\top \mathbf{x}_i \mathbf{x}_i^\top (\boldsymbol{\beta}_t - \boldsymbol{\beta}_j) \\
&= \sigma_{\max}(\mathbf{Q}_t) \sum_{i \in \mathcal{A}_i^* \cap \mathcal{A}_j} \langle \mathbf{x}_i, \boldsymbol{\beta}_t - \boldsymbol{\beta}_j \rangle^2 \\
&\leq 2\sigma_{\max}(\mathbf{Q}_t) \sum_{i \in \mathcal{A}_i^* \cap \mathcal{A}_j} (\langle \mathbf{x}_i, \boldsymbol{\beta}_t - \boldsymbol{\beta}_j^{(0)} \rangle^2 + \langle \mathbf{x}_i, \boldsymbol{\beta}_j^{(0)} - \boldsymbol{\beta}_j \rangle^2) \\
&\stackrel{(a)}{\leq} 2\sigma_{\max}(\mathbf{Q}_t) \sum_{i \in \mathcal{A}_i^* \cap \mathcal{A}_j} (\langle \mathbf{x}_i, \boldsymbol{\beta}_t - \boldsymbol{\beta}_t^{(0)} \rangle^2 + \langle \mathbf{x}_i, \boldsymbol{\beta}_j^{(0)} - \boldsymbol{\beta}_j \rangle^2) \\
&\leq 4\sigma_{\max}^2(\mathbf{Q}_t) \cdot \varepsilon_0^2,
\end{aligned}$$

where step (a) follows from the fact that for each  $i \in \mathcal{A}_i^* \cap \mathcal{A}_j$ ,  $\langle \mathbf{x}_i, \boldsymbol{\beta}_t - \boldsymbol{\beta}_j^{(0)} \rangle^2 \leq \langle \mathbf{x}_i, \boldsymbol{\beta}_t - \boldsymbol{\beta}_t^{(0)} \rangle^2$  due to the label assignment rule. Accordingly,

$$U_2 \leq 2 \sum_{t \neq j} \sigma_{\max}(\mathbf{Q}_t) \varepsilon_0. \quad (48)$$

It remains to bound  $\sigma_{\max}(\mathbf{Q}_t)$ . For each  $t$ , define

$$\mathcal{A}_j^t := \{i \in \mathcal{A}_i^* \mid |\langle \mathbf{x}_i, \boldsymbol{\beta}_t - \boldsymbol{\beta}_j^{(0)} \rangle| \leq |\langle \mathbf{x}_i, \boldsymbol{\beta}_t - \boldsymbol{\beta}_t^{(0)} \rangle|\}$$

as the set of samples that are generated from model  $t$ , but have smaller reconstruction error in  $\boldsymbol{\beta}_j^{(0)}$  compared to  $\boldsymbol{\beta}_t^{(0)}$ . We have  $\mathcal{A}_j \cap \mathcal{A}_i^* \subseteq \mathcal{A}_j^t$ , which leads to

$$\sigma_{\max}(\mathbf{Q}_t) \leq \sigma_{\max}\left(\sum_{i \in \mathcal{A}_j^t} \mathbf{x}_i \mathbf{x}_i^\top\right). \quad (49)$$

In parallel, for  $X \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ , define

$$\mathcal{E}_j^t = \{|\langle X, \boldsymbol{\beta}_t - \boldsymbol{\beta}_j^{(0)} \rangle| \leq |\langle X, \boldsymbol{\beta}_t - \boldsymbol{\beta}_t^{(0)} \rangle|\}.$$

Let  $\bar{\varepsilon}_0$  be an upper bound of  $\varepsilon_0$ .

$$\mathbb{E}[|\mathcal{A}_j^t|] = \mathbb{E}[|\mathcal{A}_i^*|] \cdot \mathbb{P}(\mathcal{E}_j^t) = \omega_t B \cdot \mathbb{P}(\mathcal{E}_j^t) \leq \omega_t B \frac{\left\| \boldsymbol{\beta}_t - \boldsymbol{\beta}_t^{(0)} \right\|_2}{\left\| \boldsymbol{\beta}_t - \boldsymbol{\beta}_j^{(0)} \right\|_2} \leq \omega_t B \frac{\bar{\varepsilon}_0}{\Delta - \bar{\varepsilon}_0} \leq \frac{2\omega_t B \bar{\varepsilon}_0}{\Delta},$$

where the first inequality follows from Lemma 6, and the last step holds when  $\bar{\varepsilon}_0 \leq \Delta/2$ . Note that  $|\mathcal{A}_j^t|$  is a summation of independent Bernoulli random variables with success probability at most  $2\omega_t \bar{\varepsilon}_0 / \Delta$ . Then by Lemma 8, we have

$$\mathbb{P}\left(|\mathcal{A}_j^t| - \mathbb{E}[|\mathcal{A}_j^t|] \geq 2\omega_t B \bar{\varepsilon}_0 / \Delta\right) \leq 2e^{-3\omega_t B \bar{\varepsilon}_0 / (4\Delta)} \leq 2e^{-3\omega_t B \bar{\varepsilon}_0 / (4\Delta)}. \quad (50)$$

Following Lemma 7 (by setting  $k = 2$ ), we have

$$\sigma_{\max}(\mathbb{E}[XX^\top \mid \mathcal{E}_j^t]) \leq 2.$$

Part (b) in Lemma 15 suggests that  $X \mid \mathcal{E}_j^t$  is still sub-Gaussian random vector with constant Orlicz norm. According to the concentration result in Remark 5.40 of [21], we have that with probability at least  $1 - 2e^{-p}$ ,

$$\sigma_{\max}\left(\sum_{i \in \mathcal{A}_j^t} \mathbf{x}_i \mathbf{x}_i^\top\right) \leq |\mathcal{A}_j^t|(2 + (\eta \vee \eta^2)),$$

where  $\eta \asymp \sqrt{p/|\mathcal{A}_j^t|}$ . We thus have

$$\sigma_{\max}\left(\sum_{i \in \mathcal{A}_j^t} \mathbf{x}_i \mathbf{x}_i^\top\right) \lesssim p \vee |\mathcal{A}_j^t| \lesssim p + |\mathcal{A}_j^t|.$$

Putting the above result, (50) and (49) together, and taking the union bound over all  $t \neq j$ , we have that with probability at least  $1 - ke^{-p} - 2ke^{-3\omega B\bar{\varepsilon}_0/(4\Delta)}$ ,

$$\sum_{t \neq j} \sigma_{\max}(\mathbf{Q}_t) \lesssim \sum_{t \neq j} p + |\mathcal{A}_j^t| \lesssim 2kp + \sum_{t \neq j} \omega_t B \bar{\varepsilon}_0 / \Delta \lesssim kp + B\bar{\varepsilon}_0 / \Delta.$$

Plugging the above result into (48) yields that for some constant  $C$

$$\mathbb{P}(U_2 \geq C(kp + B\bar{\varepsilon}_0/\Delta)\varepsilon_0) \leq ke^{-p} + 2ke^{-3\omega B\bar{\varepsilon}_0/(4\Delta)}. \quad (51)$$

**Ensemble.** Combining the bounds of  $U_1$  and  $U_2$ , there exists a constant  $C$  such that when  $B \gtrsim p/\omega$ ,

$$\mathbb{P}\left(\left\|\beta_j^{(1)} - \beta_j\right\|_2 \geq \underbrace{C \frac{(kp + B\bar{\varepsilon}_0/\Delta)\varepsilon_0}{\omega_j B}}_U\right) \leq 4e^{-C'\omega B} + 2ke^{-p} + 2ke^{-3\omega B\bar{\varepsilon}_0/(4\Delta)}.$$

Now we set  $\bar{\varepsilon}_0 = \omega\Delta/(4C)$ . Then the condition  $B \geq 4Ckp/\omega$  leads to  $U \leq \frac{1}{2}\varepsilon_0$ . Accordingly

$$\mathbb{P}\left(\left\|\beta_j^{(1)} - \beta_j\right\|_2 \geq \frac{1}{2}\varepsilon_0\right) \leq 4e^{-C'\omega B} + ke^{-p} + 2ke^{-3\omega^2 B/(16C)} \leq 2ke^{-p} + 4ke^{-C_1\omega^2 B} \leq \frac{\delta}{kT},$$

where the last step follows from conditions  $B \gtrsim \omega^{-2} \log(8k^2T/\delta)$  and  $p \geq \log(2k^2T/\delta)$ . Taking union bound over all  $j \in [k]$ , we finish proving the error decay in the first iteration. Using the same calculation for all  $T$  iterations and taking union bound concludes the proof.

## 5.6 Proof of Lemma 3

For  $X \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ , define event  $\mathcal{E}_j$ , which indicates the case that sample from model  $j$  is correctly assigned label  $j$ , as

$$\mathcal{E}_j := \left\{ |\langle X, \hat{\beta}_j - \beta_j \rangle| \leq |\langle X, \hat{\beta}_t - \beta_j \rangle|, \text{ for all } t \neq j \right\}.$$

According to (42) in the proof of Theorem 2, we have

$$\mathbb{P}(\mathcal{E}_j) \geq 1 - \sum_{t \neq j} \frac{\|\boldsymbol{\beta}_j^{(0)} - \boldsymbol{\beta}_j\|_2}{\|\boldsymbol{\beta}_t^{(0)} - \boldsymbol{\beta}_j\|_2} \geq 1 - (k-1) \frac{\widehat{\varepsilon}}{\Delta - \widehat{\varepsilon}},$$

where  $\widehat{\varepsilon} := \mathcal{E}(\{\widehat{\boldsymbol{\beta}}_j\})$ . Taking union bound over all  $n$  samples, we have that the probability of correct assignment of all labels is at least

$$1 - n(k-1) \frac{\widehat{\varepsilon}}{\Delta - \widehat{\varepsilon}} \geq 1 - \delta/2,$$

where the last step holds when  $\widehat{\varepsilon} \leq \frac{\delta}{4nk} \Delta$ . When  $n \gtrsim \frac{p}{\underline{\omega}} \vee \frac{1}{\underline{\omega}} \log(k/\delta)$ , using Lemma 8 and union bound, it is guaranteed that, with probability at least  $1 - \delta/2$ , each cluster has at least  $p$  samples. Therefore, correct label assignment will lead to exact recovery.

## 6 Proofs of Technical Lemmas

### 6.1 Proof of Lemma 5

Suppose  $\mathbf{S}$  has an SVD  $\mathbf{S} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^\top$ , where  $\mathbf{U} \in \mathbb{R}^{p \times s}$ ,  $\mathbf{V} \in \mathbb{R}^{p \times s}$  have orthonormal columns  $\mathbf{U}^\top \mathbf{U} = \mathbf{V}^\top \mathbf{V} = \mathbf{I}_s$ . We can always find  $\boldsymbol{\alpha} \in \mathbb{R}^s$ ,  $\boldsymbol{\gamma} \in \mathbb{R}^p$  such that  $\boldsymbol{\beta} = \mathbf{U}\boldsymbol{\alpha} + \boldsymbol{\gamma}$ , where  $\mathbf{U}^\top \boldsymbol{\gamma} = \mathbf{0}$  and  $\|\boldsymbol{\alpha}\|_2^2 + \|\boldsymbol{\gamma}\|_2^2 = 1$ . We let  $X \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ ,  $Y = \langle X, \boldsymbol{\beta} \rangle$ .

**Proof of Inequality (27).** We find

$$\begin{aligned} \left\| \mathbf{S}^\top (\mathbf{m}_1 - \overline{\mathbf{m}}_1) \right\|_2 &= \left\| \mathbf{V}\boldsymbol{\Sigma}\mathbf{U}^\top \left( \frac{1}{n} \sum_{i=1}^n y_i^3 \mathbf{x}_i - \mathbb{E}[Y^3 X] \right) \right\|_2 \\ &= \left\| \boldsymbol{\Sigma} \left( \frac{1}{n} \sum_{i=1}^n y_i^3 \mathbf{U}^\top \mathbf{x}_i - \mathbb{E}[Y^3 \mathbf{U}^\top X] \right) \right\|_2 \\ &\leq \|\mathbf{S}\|_{op} \cdot \left\| \frac{1}{n} \sum_{i=1}^n y_i^3 \mathbf{U}^\top \mathbf{x}_i - \mathbb{E}[Y^3 \mathbf{U}^\top X] \right\|_2 \\ &= \|\mathbf{S}\|_{op} \cdot \left\| \frac{1}{n} \sum_{i=1}^n (\langle \boldsymbol{\alpha}, \tilde{\mathbf{x}}_i \rangle + z_i)^3 \tilde{\mathbf{x}}_i - \mathbb{E}[(\langle \boldsymbol{\alpha}, \tilde{X} \rangle + Z)^3 \tilde{X}] \right\|_2, \end{aligned} \quad (52)$$

where we let  $\tilde{X} := \mathbf{U}^\top X$ ,  $Z := \langle \boldsymbol{\gamma}, X \rangle$ , and  $\{(z_i, \tilde{\mathbf{x}}_i)\}_{i=1}^n$  are  $n$  independent samples of  $(\tilde{X}, Z)$ . Thanks to the rotation invariance of Gaussian, we have  $\tilde{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_s)$  and  $Z \sim \mathcal{N}(0, \|\boldsymbol{\gamma}\|_2^2)$ . Moreover,  $\tilde{X}$  and  $Z$  are independent since  $\mathbf{U}^\top \boldsymbol{\gamma} = \mathbf{0}$ .

For any  $\tau_1, \tau_2 > 1$ , define events

$$\mathcal{E} := \left\{ |\langle \boldsymbol{\alpha}, \tilde{X} \rangle| \leq \tau_1, |Z| \leq \tau_2 \right\}, \quad \mathcal{E}_n := \left\{ |\langle \boldsymbol{\alpha}, \tilde{\mathbf{x}}_i \rangle| \leq \tau_1, |z_i| \leq \tau_2, \text{ for all } i \in [n] \right\}. \quad (53)$$

We have

$$\begin{aligned}
& \left\| \frac{1}{n} \sum_{i=1}^n (\langle \boldsymbol{\alpha}, \tilde{\mathbf{x}}_i \rangle + z_i)^3 \tilde{\mathbf{x}}_i - \mathbb{E} \left[ (\langle \boldsymbol{\alpha}, \tilde{X} \rangle + Z)^3 \tilde{X} \right] \right\|_2 \\
& \leq \underbrace{\left\| \frac{1}{n} \sum_{i=1}^n (\langle \boldsymbol{\alpha}, \tilde{\mathbf{x}}_i \rangle + z_i)^3 \tilde{\mathbf{x}}_i - \mathbb{E} \left[ (\langle \boldsymbol{\alpha}, \tilde{X} \rangle + Z)^3 \tilde{X} \mid \mathcal{E} \right] \right\|_2}_{d_1} \\
& \quad + \underbrace{\left\| \mathbb{E} \left[ (\langle \boldsymbol{\alpha}, \tilde{X} \rangle + Z)^3 \tilde{X} \mid \mathcal{E} \right] - \mathbb{E} \left[ (\langle \boldsymbol{\alpha}, \tilde{X} \rangle + Z)^3 \tilde{X} \right] \right\|_2}_{d_2}.
\end{aligned}$$

For term  $d_2$ , using (62) in Lemma 13 by replacing  $(a, b, \tau_1, \tau_2)$  in the statement with

$$(\|\boldsymbol{\alpha}\|_2, \|\boldsymbol{\gamma}\|_2, \tau_1/\|\boldsymbol{\alpha}\|_2, \tau_2/\|\boldsymbol{\gamma}\|_2),$$

we obtain

$$d_2 \leq \tau_1 \left( \frac{\tau_1}{\|\boldsymbol{\alpha}\|_2} e^{-\frac{\tau_1^2}{2\|\boldsymbol{\alpha}\|_2^2}} + \frac{\tau_2}{\|\boldsymbol{\gamma}\|_2} e^{-\frac{\tau_2^2}{2\|\boldsymbol{\gamma}\|_2^2}} \right) \leq \tau_1 (\tau_1 e^{-\tau_1^2/2} + \tau_2 e^{-\tau_2^2/2}),$$

where the last step follows from the fact that function  $xe^{-x^2/2}$  is monotonically decreasing on  $x \geq 1$ . To ease notation, we let

$$\tilde{X}' \sim X \mid \mathcal{E}, \quad Z' \sim Z \mid \mathcal{E}. \quad (54)$$

Suppose  $\{(\tilde{\mathbf{x}}'_i, z'_i)\}_{i=1}^n$  are independent samples of  $(\tilde{X}', Z')$ . We observe that

$$\mathbb{P}(d_1 \geq t) \leq \mathbb{P} \left( \left\| \frac{1}{n} \sum_{i=1}^n (\langle \boldsymbol{\alpha}, \tilde{\mathbf{x}}'_i \rangle + z'_i)^3 \tilde{\mathbf{x}}'_i - \mathbb{E} \left[ (\langle \boldsymbol{\alpha}, \tilde{X}' \rangle + Z')^3 \tilde{X}' \right] \right\|_2 \geq t \right) + \mathbb{P}(\mathcal{E}_n^c).$$

Since  $|\langle \boldsymbol{\alpha}, X' \rangle + Z'| \leq \tau_1 + \tau_2$ ,  $(\langle \boldsymbol{\alpha}, \tilde{X}' \rangle + Z')^3 \tilde{X}'$  is sub-Gaussian random vector with Orlicz norm

$$\left\| (\langle \boldsymbol{\alpha}, \tilde{X}' \rangle + Z')^3 \tilde{X}' \right\|_{\psi_2} \lesssim (\tau_1 + \tau_2)^3.$$

By concentration result (58) in Lemma 10, we have that for some constants  $C_1, C_2$ , condition  $n \geq C_1 s (\tau_1 + \tau_2)^6 / t^2$  leads to

$$\mathbb{P} \left( \left\| \frac{1}{n} \sum_{i=1}^n (\langle \boldsymbol{\alpha}, \tilde{\mathbf{x}}'_i \rangle + z'_i)^3 \tilde{\mathbf{x}}'_i - \mathbb{E} \left[ (\langle \boldsymbol{\alpha}, \tilde{X}' \rangle + Z')^3 \tilde{X}' \right] \right\|_2 \geq t \right) \leq e^{-C_2 n t^2 / (\tau_1 + \tau_2)^6}.$$

Meanwhile, the variance of  $\langle \boldsymbol{\alpha}, \tilde{X} \rangle$  and  $Z$  are both at most 1. We thus obtain

$$\mathbb{P}(\mathcal{E}_n^c) \leq n e^{1-\tau_1^2} + n e^{1-\tau_2^2} \quad (55)$$

by using Gaussian tail bound and union bound. Accordingly,

$$\mathbb{P}(d_1 \geq t) \leq e^{-c_2 n t^2 / (\tau_1 + \tau_2)^6} + n e^{1-\tau_1^2} + n e^{1-\tau_2^2}.$$

Setting  $\tau_1 = \tau_2 = C\sqrt{\log n}$  for sufficiently large constant  $C$  and  $t \asymp (\tau_1 + \tau_2)^3 \sqrt{1/n} \left( \sqrt{\log(\frac{2}{\delta})} \vee \sqrt{s} \right)$ , we have  $d_2 \lesssim 1/n$  and  $\mathbb{P}(d_1 \geq t) \leq \delta/2 + 1/n$ . Requiring  $n \gtrsim 2/\delta$  gives our result.

**Proof of Inequality (28).** We find

$$\begin{aligned} \left\| \mathbf{S}^\top (\mathbf{M}_2 - \overline{\mathbf{M}}_2) \mathbf{S} \right\|_{op} &= \left\| \mathbf{V} \boldsymbol{\Sigma} \mathbf{U}^\top \left( \frac{1}{n} \sum_{i \in [n]} y_i^2 \mathbf{x}_i \mathbf{x}_i^\top - \mathbb{E} \left[ Y^2 \mathbf{X} \mathbf{X}^\top \right] \right) \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\top \right\|_{op} \\ &\leq \|\mathbf{S}\|_{op}^2 \cdot \left\| \frac{1}{n} \sum_{i \in [n]} y_i^2 \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top - \mathbb{E} \left[ Y^2 \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top \right] \right\|_{op}, \end{aligned}$$

where  $\tilde{\mathbf{x}}_i$  and  $\tilde{\mathbf{X}}$  are defined according to (52). Using the  $\mathcal{E}, \mathcal{E}_n$  defined in (53), we have

$$\begin{aligned} &\left\| \frac{1}{n} \sum_{i \in [n]} y_i^2 \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top - \mathbb{E} \left[ Y^2 \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top \right] \right\|_{op} \\ &\leq \underbrace{\left\| \frac{1}{n} \sum_{i=1}^n (\langle \boldsymbol{\alpha}, \tilde{\mathbf{x}}_i \rangle + z_i)^2 \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top - \mathbb{E} \left[ (\langle \boldsymbol{\alpha}, \tilde{\mathbf{X}} \rangle + Z)^2 \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top \mid \mathcal{E} \right] \right\|_{op}}_{d_1} \\ &\quad + \underbrace{\left\| \mathbb{E} \left[ (\langle \boldsymbol{\alpha}, \tilde{\mathbf{X}} \rangle + Z)^2 \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top \mid \mathcal{E} \right] - \mathbb{E} \left[ (\langle \boldsymbol{\alpha}, \tilde{\mathbf{X}} \rangle + Z)^2 \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top \right] \right\|_{op}}_{d_2}. \end{aligned}$$

Applying (63) in Lemma 13 via setting  $(a, b, \tau_1, \tau_2)$  in the statement to be

$$(\|\boldsymbol{\alpha}\|_2, \|\boldsymbol{\gamma}\|_2, \tau_1/\|\boldsymbol{\alpha}\|_2, \tau_2/\|\boldsymbol{\gamma}\|_2)$$

provides that

$$d_2 \leq \frac{\tau_1^3}{\|\boldsymbol{\alpha}\|_2^3} e^{-\frac{\tau_1^2}{2\|\boldsymbol{\alpha}\|_2^2}} + \frac{\tau_1}{\|\boldsymbol{\alpha}\|_2} \frac{\tau_2}{\|\boldsymbol{\gamma}\|_2} e^{-\frac{\tau_1^2}{2\|\boldsymbol{\alpha}\|_2^2} - \frac{\tau_2^2}{2\|\boldsymbol{\gamma}\|_2^2}} \leq \tau_1^3 e^{-\tau_1^2/2} + \tau_1 \tau_2 e^{-\tau_1^2/2 - \tau_2^2/2},$$

where the last inequality follows from the fact that functions  $x^3 e^{-x^2/2}$ ,  $x e^{-x^2/2}$  are monotonically decreasing when  $x$  is sufficiently large.

We follow the same idea used before to bound  $d_1$ . Introduce  $\tilde{\mathbf{X}}', Z'$  according to (54). Then we obtain

$$\mathbb{P}(d_1 \geq t) \leq \mathbb{P} \left( \left\| \frac{1}{n} \sum_{i=1}^n (\langle \boldsymbol{\alpha}, \tilde{\mathbf{x}}'_i \rangle + z'_i)^2 \tilde{\mathbf{x}}'_i \tilde{\mathbf{x}}'^{\top}_i - \mathbb{E} \left[ (\langle \boldsymbol{\alpha}, \tilde{\mathbf{X}}' \rangle + Z')^2 \tilde{\mathbf{X}}' \tilde{\mathbf{X}}'^{\top} \right] \right\|_{op} \geq t \right) + \mathbb{P}(\mathcal{E}_n^c). \quad (56)$$

Since  $|\langle \boldsymbol{\alpha}, \tilde{\mathbf{X}}' \rangle + Z'| \leq \tau_1 + \tau_2$ ,  $(\langle \boldsymbol{\alpha}, \tilde{\mathbf{X}}' \rangle + Z') \tilde{\mathbf{X}}'$  is sub-Gaussian random vector with norm  $\mathcal{O}(\tau_1 + \tau_2)$ . Applying (58) in Lemma 10, we have that for  $t \in (0, (\tau_1 + \tau_2)^2)$  and some constants  $C_1, C_2$ , the condition  $n \geq C_1 k (\tau_1 + \tau_2)^4 / t^2$  yields

$$\mathbb{P} \left( \left\| \frac{1}{n} \sum_{i=1}^n (\langle \boldsymbol{\alpha}, \tilde{\mathbf{x}}'_i \rangle + z'_i)^2 \tilde{\mathbf{x}}'_i \tilde{\mathbf{x}}'^{\top}_i - \mathbb{E} \left[ (\langle \boldsymbol{\alpha}, \tilde{\mathbf{X}}' \rangle + Z')^2 \tilde{\mathbf{X}}' \tilde{\mathbf{X}}'^{\top} \right] \right\|_{op} \geq t \right) \leq e^{-C_2 n t^2 / (\tau_1 + \tau_2)^4}.$$



Plugging it back into (56) and using the bound (55) of  $\mathbb{P}(\mathcal{E}^c)$ , we obtain

$$\mathbb{P}(d_1 \geq t) \leq e^{-c_2 n t^2 / (\tau_1 + \tau_2)^4} + n e^{1 - \tau_1^2} + n e^{1 - \tau_2^2}.$$

Choosing  $\tau_1 = \tau_2 = C\sqrt{\log n}$  for sufficiently large constant  $C$  and letting  $t \asymp \frac{\log n}{\sqrt{n}} \left( \sqrt{\log\left(\frac{2}{\delta}\right)} \vee \sqrt{s} \right)$ , we have that when  $n \geq C'(1/\delta \vee s)$  for sufficiently large  $C'$ , it is guaranteed that  $\mathbb{P}(d_1 \geq t) \leq \delta$  and  $d_2 \lesssim 1/n$ , which concludes the proof.

**Proof of Inequality (29).** Using the Cauchy-Schwartz inequality and the definitions of  $\tilde{\mathbf{x}}_i$  and  $\tilde{X}$  in (52), we have that

$$\|(\mathbf{M}_3 - \overline{\mathbf{M}}_3)(\mathbf{S}, \mathbf{S}, \mathbf{S})\|_{op} \leq \|\mathbf{S}\|_{op}^3 \cdot \left\| \frac{1}{n} \sum_{i \in [n]} y_i^3 \tilde{\mathbf{x}}_i \otimes \tilde{\mathbf{x}}_i \otimes \tilde{\mathbf{x}}_i - \mathbb{E} \left[ Y^3 \tilde{X} \otimes \tilde{X} \otimes \tilde{X} \right] \right\|_{op}.$$

Again, we use the event  $\mathcal{E}, \mathcal{E}_n$  in (53) to bound the operator norm. In detail, we have

$$\begin{aligned} \left\| \frac{1}{n} \sum_{i \in [n]} y_i^3 \tilde{\mathbf{x}}_i^{\otimes 3} - \mathbb{E} \left[ Y^3 \tilde{X}^{\otimes 3} \right] \right\|_{op} &\leq \underbrace{\left\| \frac{1}{n} \sum_{i=1}^n (\langle \boldsymbol{\alpha}, \tilde{\mathbf{x}}_i \rangle + z_i)^3 \tilde{\mathbf{x}}_i^{\otimes 3} - \mathbb{E} \left[ (\langle \boldsymbol{\alpha}, \tilde{X} \rangle + Z)^3 \tilde{X}^{\otimes 3} \mid \mathcal{E} \right] \right\|_{op}}_{d_1} \\ &+ \underbrace{\left\| \mathbb{E} \left[ (\langle \boldsymbol{\alpha}, \tilde{X} \rangle + Z)^3 \tilde{X}^{\otimes 3} \mid \mathcal{E} \right] - \mathbb{E} \left[ (\langle \boldsymbol{\alpha}, \tilde{X} \rangle + Z)^3 \tilde{X}^{\otimes 3} \right] \right\|_{op}}_{d_2}. \end{aligned}$$

Applying (64) in Lemma 13 by setting  $(a, b, \tau_1, \tau_2)$  in the statement to be  $(\|\boldsymbol{\alpha}\|_2, \|\boldsymbol{\gamma}\|_2, \tau_1/\|\boldsymbol{\alpha}\|_2, \tau_2/\|\boldsymbol{\gamma}\|_2)$ , we obtain

$$d_2 \leq \frac{\tau_1^5}{\|\boldsymbol{\alpha}\|_2^5} e^{-\frac{\tau_1^2}{2\|\boldsymbol{\alpha}\|_2^2}} + \frac{\tau_1^3}{\|\boldsymbol{\alpha}\|_2^3} \frac{\tau_2}{\|\boldsymbol{\gamma}\|_2} e^{-\frac{\tau_1^2}{2\|\boldsymbol{\gamma}\|_2^2} - \frac{\tau_2^2}{2\|\boldsymbol{\gamma}\|_2^2}} \leq \tau_1^5 e^{-\tau_1^2/2} + \tau_1^3 \tau_2 e^{-\tau_1^2/2 - \tau_2^2/2},$$

where the last inequality follows from the fact that functions  $x^5 e^{-x^2/2}$ ,  $x^3 e^{-x^2/2}$ ,  $x e^{-x^2/2}$  are monotonically decreasing when  $x$  is sufficiently large. For term  $d_1$ , introducing the  $\tilde{X}'$ ,  $Z'$  according to (54), we have

$$\mathbb{P}(d_1 \geq t) \leq \mathbb{P} \left( \left\| \frac{1}{n} \sum_{i=1}^n (\langle \boldsymbol{\alpha}, \tilde{\mathbf{x}}_i \rangle + z_i)^3 \tilde{\mathbf{x}}_i'^{\otimes 3} - \mathbb{E} \left[ (\langle \boldsymbol{\alpha}, \tilde{X}' \rangle + Z')^3 \tilde{X}'^{\otimes 3} \right] \right\|_{op} \geq t \right) + \mathbb{P}(\mathcal{E}_n^c).$$

Note that  $(\langle \boldsymbol{\alpha}, \tilde{X}' \rangle + Z')\tilde{X}'$  is sub-Gaussian random vector with norm  $\mathcal{O}(\tau_1 + \tau_2)$ . Applying (60) in Lemma 10, we find that for any  $t \in (0, (\tau_1 + \tau_2)^3 \sqrt{s})$  and constants  $C_1, C_2$ , condition  $n \geq C_1(\tau_1 + \tau_2)^6 s^2/t^2$  yields

$$\mathbb{P}(d_1 \geq t) \leq e^{-c_2 \frac{nt^2}{k^2(\tau_1 + \tau_2)^6}} + \mathbb{P}(\mathcal{E}_n^c) \leq e^{-c_2 \frac{nt^2}{k^2(\tau_1 + \tau_2)^6}} + n e^{1 - \tau_1^2} + n e^{1 - \tau_2^2}.$$

Setting  $\tau_1 = \tau_2 = C\sqrt{\log n}$  for sufficiently large constant  $C$ ,  $t \asymp \frac{s\sqrt{\log n}^3}{\sqrt{n}} \sqrt{\log\left(\frac{2}{\delta}\right)}$ , and assuming  $n \gtrsim \max\{s \log\left(\frac{2}{\delta}\right), 1/\delta\}$ , we obtain that  $\mathbb{P}(d_1 \leq t) \leq \delta$  and  $d_2 \leq 1/n$ , which concludes the proof.

## 6.2 Proof of Lemma 6

For two vector  $\mathbf{u}, \mathbf{v}$ , we define angle  $\alpha(\mathbf{u}, \mathbf{v}) \in [0, \pi]$  as

$$\alpha(\mathbf{u}, \mathbf{v}) := \cos^{-1} \frac{(\mathbf{u} - \mathbf{v})^\top (\mathbf{u} + \mathbf{v})}{\|\mathbf{u} + \mathbf{v}\|_2 \cdot \|\mathbf{u} - \mathbf{v}\|_2}.$$

Without loss of generality, we assume  $\mathbf{u}, \mathbf{v}$  live in the subspace spanned by  $\mathbf{e}_1, \mathbf{e}_2$ . We use  $x_1, x_2$  to denote the first two coordinates of  $X$ . We can let

$$x_1 = A \cos \theta, \quad x_2 = A \sin \theta,$$

where  $A$  is Rayleigh random variable, and  $\theta$  is uniformly distributed over  $[0, 2\pi)$ . Conditioning on  $\mathcal{E}$ , the range of  $\theta$  is truncated to be  $[\theta_0, \theta_0 + \alpha(\mathbf{u}, \mathbf{v})] \cup [\theta_0 + \pi, \theta_0 + \pi + \alpha(\mathbf{u}, \mathbf{v})]$ , where  $\theta_0$  depends on  $\mathbf{u}, \mathbf{v}$ . Therefore, we have

$$\mathbb{P}(\mathcal{E}) = \frac{\alpha(\mathbf{u}, \mathbf{v})}{\pi}.$$

If  $\|\mathbf{u}\|_2 > \|\mathbf{v}\|_2$ ,

$$\cos[\alpha(\mathbf{u}, \mathbf{v})] \geq \frac{\|\mathbf{u}\|_2^2 - \|\mathbf{v}\|_2^2}{\|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2} > 0.$$

So we have  $\alpha(\mathbf{u}, \mathbf{v}) \in [0, \pi/2]$ . Using the fact that  $\alpha < \frac{\pi}{2} \sin \alpha$  for any  $\alpha \in [0, \pi/2]$ , we have

$$\mathbb{P}(\mathcal{E}) \leq \frac{1}{2} \sin[\alpha(\mathbf{u}, \mathbf{v})] \leq \frac{\|\mathbf{u}\|_2 \|\mathbf{v}\|_2}{\|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2} \leq \frac{\|\mathbf{v}\|_2}{\|\mathbf{u}\|_2}.$$

## 6.3 Proof of Lemma 7

Note that conditioning on  $\mathcal{E}$  or  $\mathcal{E}^c$  will not change the distribution of  $\|X\|_2$ . We thus have

$$\mathbb{E} \left[ \|X\|_2^2 \mid \mathcal{E} \right] = \mathbb{E} \left[ \|X\|_2^2 \mid \mathcal{E}^c \right] = \mathbb{E} \|X\|_2^2 = p.$$

Hence,

$$\text{Trace} \left( \mathbb{E} \left[ XX^\top \mid \mathcal{E} \right] \right) = p. \tag{57}$$

Also note that  $\mathbb{E} [XX^\top \mid \mathcal{E}]$  and  $\mathbb{E} [XX^\top \mid \mathcal{E}^c]$  have at least  $p - k$  eigenvalues that are 1 since  $\{\mathbf{u}_1, \dots, \mathbf{u}_k\}$  spans a subspace with dimension at most  $k$ . Therefore we have

$$\sigma_{max} \left( \mathbb{E} \left[ XX^\top \mid \mathcal{E} \right] \right) \leq \text{Trace} \left( \mathbb{E} \left[ XX^\top \mid \mathcal{E} \right] \right) - (p - k) \leq k.$$

The above inequality also holds for  $\sigma_{max} (\mathbb{E} [XX^\top \mid \mathcal{E}^c])$ . Note that

$$\mathbf{I}_p = \mathbb{E}[XX^\top] = \mathbb{E} \left[ XX^\top \mid \mathcal{E} \right] \mathbb{P}(\mathcal{E}) + \mathbb{E} \left[ XX^\top \mid \mathcal{E}^c \right] (1 - \mathbb{P}(\mathcal{E})).$$

Suppose  $\mathbf{v}$  is the eigenvector that corresponds to the minimum eigenvalue of  $\mathbb{E} [XX^\top \mid \mathcal{E}]$ . Therefore, we have

$$\begin{aligned} 1 &= \mathbf{v}^\top \mathbb{E} \left[ XX^\top \mid \mathcal{E} \right] \mathbf{v} \mathbb{P}(\mathcal{E}) + \mathbf{v}^\top \mathbb{E} \left[ XX^\top \mid \mathcal{E}^c \right] \mathbf{v} (1 - \mathbb{P}(\mathcal{E})) \\ &\leq \sigma_{min} \left( \mathbb{E} \left[ XX^\top \mid \mathcal{E} \right] \right) \mathbb{P}(\mathcal{E}) + \mathbf{v}^\top \mathbb{E} \left[ XX^\top \mid \mathcal{E}^c \right] \mathbf{v} (1 - \mathbb{P}(\mathcal{E})) \\ &\leq \sigma_{min} \left( \mathbb{E} \left[ XX^\top \mid \mathcal{E} \right] \right) \mathbb{P}(\mathcal{E}) + k(1 - \mathbb{P}(\mathcal{E})). \end{aligned}$$

## 7 Auxiliary Results

**Lemma 8** (Sum of Bernoulli Random Variables). *Suppose  $X_1, \dots, X_n$  are  $n$  independent Bernoulli random variables with  $\mathbb{P}[X_1 = 0] = 1 - p$  and  $\mathbb{P}[X_1 = 1] = p$ . Let*

$$\bar{X} = \frac{1}{n} \sum_{i \in [n]} X_i.$$

For every  $t > 0$ , we have

$$\mathbb{P}(|\bar{X} - \mathbb{E}[\bar{X}]| \geq tp) \leq 2e^{-\frac{3t^2}{2(t+3)}np}.$$

*Proof.* We find that  $X_1 - \mathbb{E}[X_1]$  has variance  $p(1-p)$  and  $|X_1 - \mathbb{E}[X_1]| \leq 1$ . Using Bernstein's inequality, we have

$$\mathbb{P}(|\bar{X} - \mathbb{E}[\bar{X}]| \geq tp) \leq 2e^{-\frac{nt^2p^2/2}{p(1-p)+tp/3}} \leq 2e^{-\frac{3t^2}{2(t+3)}np}.$$

□

**Lemma 9** (Properties of Whitening Matrices, Lemma 6 in [6]). *Suppose  $\mathbf{A}$  and  $\widehat{\mathbf{A}}$  are both positive semidefinite matrices in  $\mathbb{R}^{p \times p}$  with rank  $k$ . Let  $\mathbf{W}, \widehat{\mathbf{W}} \in \mathbb{R}^{p \times k}$  be whitening matrices such that  $\mathbf{W}^\top \mathbf{A} \mathbf{W} = \mathbf{I}_k$ ,  $\widehat{\mathbf{W}}^\top \widehat{\mathbf{A}} \widehat{\mathbf{W}} = \mathbf{I}_k$ . When  $\alpha := \left\| \mathbf{A} - \widehat{\mathbf{A}} \right\|_{op} / \sigma_k(\mathbf{A}) < 1/3$ , we have*

$$\begin{aligned} \left\| \widehat{\mathbf{W}} \right\|_{op} &\leq 2 \left\| \mathbf{W} \right\|_{op}, & \left\| \widehat{\mathbf{W}}^\dagger \right\|_{op} &\leq 2 \left\| \mathbf{W}^\dagger \right\|_{op}, \\ \left\| \mathbf{W} - \widehat{\mathbf{W}} \right\|_{op} &\leq 2\alpha \cdot \left\| \mathbf{W} \right\|_{op}, & \left\| \mathbf{W}^\dagger - \widehat{\mathbf{W}}^\dagger \right\|_{op} &\leq 2\alpha \cdot \left\| \mathbf{W}^\dagger \right\|_{op}. \end{aligned}$$

**Lemma 10** (Concentration of Sub-Gaussian Vectors). *Suppose  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^p$  are  $n$  i.i.d. sub-Gaussian vectors with Orlicz norm  $\|\mathbf{x}_1\|_{\psi_2} \leq K$ .*

1. *There exist constants  $C_i$  such that for every  $t > 0$ , when  $n \geq C_1(K/t)^2p$ ,*

$$\mathbb{P} \left( \left\| \frac{1}{n} \sum_{i \in [n]} \mathbf{x}_i - \mathbb{E}[\mathbf{x}_1] \right\|_2 \geq t \right) \leq e^{-C_2nt^2/K^2}. \quad (58)$$

2. *There exist constants  $C_i$  such that for every  $t \in (0, K^2)$ , when  $n \geq C_1(K^2/t)^2p$ ,*

$$\mathbb{P} \left( \left\| \frac{1}{n} \sum_{i \in [n]} \mathbf{x}_i \mathbf{x}_i^\top - \mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top] \right\|_{op} \geq t \right) \leq e^{-C_2nt^2/K^4}. \quad (59)$$

3. *There exist constants  $C_i$  such that for every  $t \in (0, K^3\sqrt{p})$ , when  $n \geq C_1(K^3/t)^2p^2$ ,*

$$\mathbb{P} \left( \left\| \frac{1}{n} \sum_{i \in [n]} \mathbf{x}_i^{\otimes 3} - \mathbb{E}[\mathbf{x}_1^{\otimes 3}] \right\|_{op} \geq t \right) \leq e^{-C_2nt^2/(p^2K^6)}. \quad (60)$$

*Proof.*

1. Note that

$$\left\| \frac{1}{n} \sum_{i \in [n]} \mathbf{x}_i - \mathbb{E}[\mathbf{x}_1] \right\|_2 = \sup_{\mathbf{u} \in \mathbb{S}^{p-1}} \left| \frac{1}{n} \sum_{i=1}^n \langle \mathbf{x}_i, \mathbf{u} \rangle - \mathbb{E}[\langle \mathbf{x}_i, \mathbf{u} \rangle] \right|.$$

Since  $\mathbf{x}_i$  is sub-Gaussian vector, then for any fixed  $\mathbf{u} \in \mathbb{S}^{p-1}$ ,  $\langle \mathbf{x}_i, \mathbf{u} \rangle$  is sub-Gaussian random variable with norm  $K$ . Therefore,  $\langle \mathbf{x}_i, \mathbf{u} \rangle - \mathbb{E}[\langle \mathbf{x}_i, \mathbf{u} \rangle]$  is also sub-Gaussian with norm at most  $2K$ . By standard concentration of sub-Gaussianity, for some constant  $C$ , we obtain

$$\mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n \langle \mathbf{x}_i, \mathbf{u} \rangle - \mathbb{E}[\langle \mathbf{x}_i, \mathbf{u} \rangle] \right| \geq t \right) \leq e^{1-Cnt^2/K^2}.$$

It is possible to construct an  $\epsilon$ -net  $\mathcal{S}_\epsilon$  of  $\mathbb{S}^{p-1}$  with size  $|\mathcal{S}_\epsilon| \leq (1 + 2/\epsilon)^p$  (Lemma 5.2 in [21]). Applying probabilistic union bound leads to

$$\mathbb{P} \left( \sup_{\mathbf{u} \in \mathcal{S}_\epsilon} \left| \frac{1}{n} \sum_{i=1}^n \langle \mathbf{x}_i, \mathbf{u} \rangle - \mathbb{E}[\langle \mathbf{x}_i, \mathbf{u} \rangle] \right| \geq t \right) \leq (1 + 2/\epsilon)^p e^{1-Cnt^2/K^2}.$$

For any  $\mathbf{z} \in \mathbb{S}^{p-1}$ , we can always find  $\mathbf{u} \in \mathcal{S}_\epsilon$  such that  $\|\mathbf{z} - \mathbf{u}\|_2 \leq \epsilon$ . Then

$$\left| \frac{1}{n} \sum_{i=1}^n \langle \mathbf{x}_i, \mathbf{z} \rangle - \mathbb{E}[\langle \mathbf{x}_i, \mathbf{z} \rangle] \right| \leq \left| \frac{1}{n} \sum_{i=1}^n \langle \mathbf{x}_i, \mathbf{u} \rangle - \mathbb{E}[\langle \mathbf{x}_i, \mathbf{u} \rangle] \right| + \|\mathbf{u} - \mathbf{z}\|_2 \cdot \left\| \frac{1}{n} \sum_{i \in [n]} \mathbf{x}_i - \mathbb{E}[\mathbf{x}_1] \right\|_2.$$

Therefore, we obtain

$$\sup_{\mathbf{z} \in \mathbb{S}^{p-1}} \left| \frac{1}{n} \sum_{i=1}^n \langle \mathbf{x}_i, \mathbf{z} \rangle - \mathbb{E}[\langle \mathbf{x}_i, \mathbf{z} \rangle] \right| \leq \frac{1}{1-\epsilon} \cdot \sup_{\mathbf{u} \in \mathcal{S}_\epsilon} \left| \frac{1}{n} \sum_{i=1}^n \langle \mathbf{x}_i, \mathbf{u} \rangle - \mathbb{E}[\langle \mathbf{x}_i, \mathbf{u} \rangle] \right|. \quad (61)$$

Setting  $\epsilon = 1/4$  and assuming  $n \geq C'(K/t)^2 p$  for sufficiently large constant  $C'$  completes the proof.

2. Refer to Theorem 5.39 in [21] for the proof.

3. Note that for any 3-way tensor  $\mathbf{T} \in \mathbb{R}^{p \times p \times p}$  and two vectors  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^p$  that satisfy  $\|\mathbf{u} - \mathbf{v}\|_2 \leq \epsilon$ , we have

$$\begin{aligned} \mathbf{T}(\mathbf{u}, \mathbf{u}, \mathbf{u}) - \mathbf{T}(\mathbf{v}, \mathbf{v}, \mathbf{v}) &= \mathbf{T}(\mathbf{u} - \mathbf{v}, \mathbf{u}, \mathbf{u}) + \mathbf{T}(\mathbf{v}, \mathbf{u} - \mathbf{v}, \mathbf{u}) + \mathbf{T}(\mathbf{v}, \mathbf{v}, \mathbf{u} - \mathbf{v}) \\ &\leq 3\epsilon \cdot \sup_{\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathbb{S}^{p-1}} |\mathbf{T}(\mathbf{a}, \mathbf{b}, \mathbf{c})| \leq 27\epsilon \|\mathbf{T}\|_{op}, \end{aligned}$$

where the last inequality follows from Lemma 12. Constructing an  $\epsilon$ -net  $\mathcal{S}_\epsilon$  on  $\mathbb{S}^{p-1}$  and following similar idea in showing (61), we obtain

$$\left\| \frac{1}{n} \sum_{i \in [n]} \mathbf{x}_i^{\otimes 3} - \mathbb{E}[\mathbf{x}_1^{\otimes 3}] \right\|_{op} \leq \frac{1}{1-27\epsilon} \sup_{\mathbf{u} \in \mathcal{S}_\epsilon} \left| \frac{1}{n} \sum_{i \in [n]} \langle \mathbf{x}_i, \mathbf{u} \rangle^3 - \mathbb{E}[\langle \mathbf{x}_1, \mathbf{u} \rangle^3] \right|.$$

Now we set  $\epsilon = 1/54$ , which leads to  $|\mathcal{S}_\epsilon| \leq 109^p$ . For any fixed  $\mathbf{u} \in \mathbb{S}^{p-1}$ ,  $\langle \mathbf{x}_i, \mathbf{u} \rangle$  is sub-Gaussian random variable with norm  $K$ . Using the concentration of cubes of sub-Gaussians (Lemma 11) and applying union bound, we obtain

$$\mathbb{P} \left( \sup_{\mathbf{u} \in \mathcal{S}_\epsilon} \left| \frac{1}{n} \sum_{i \in [n]} \langle \mathbf{x}_i, \mathbf{u} \rangle^3 - \mathbb{E} [\langle \mathbf{x}_1, \mathbf{u} \rangle^3] \right| > CK^3 \frac{\sqrt{p^3 \log^3(109/\delta) + 2p^2 \log^2(109/\delta)n}}{n} \right) \leq \delta$$

for any  $\delta \in (0, 1)$  and some constant  $C > 0$ . Finally, for any  $t \in (0, K^3 \sqrt{p})$ , setting  $\delta = e^{-C' \frac{nt^2}{p^2 K^6}}$ ,  $n \geq C''(p/t)^2 K^6$  for some constants  $C', C''$  completes the proof.  $\square$

The next result shows a tail bound of a finite sum of sub-Gaussian random variables. A similar result is proved in the case of Gaussian in [12]. Here, we present our proof that can cover general sub-Gaussian distribution.

**Lemma 11** (Sum of Cubes of Sub-Gaussians). *Suppose  $X_1, X_2, \dots, X_n$  are  $n$  i.i.d. sub-Gaussian random variables with Orlicz norm  $\|X_1\|_{\psi_2} \leq K$ . There exists an absolute constant  $C$  such that for any  $\delta \in (0, 1)$ ,*

$$\mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n X_i^3 - \mathbb{E} [X_1^3] \right| > CK^3 \frac{\sqrt{\log^3(1/\delta) + 2 \log^2(1/\delta)n}}{n} \right) \leq \delta.$$

*Proof.* For any positive even integer  $q$  and  $t \in \mathbb{R}^+$ , by Markov's inequality, we have

$$\begin{aligned} \mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n X_i^3 - \mathbb{E} [X_1^3] \right| > t \right) &= \mathbb{P} \left( \left( \frac{1}{n} \sum_{i=1}^n X_i^3 - \mathbb{E} [X_1^3] \right)^q > t^q \right) \\ &\leq \frac{1}{t^q} \mathbb{E} \left[ \left( \frac{1}{n} \sum_{i=1}^n X_i^3 - \mathbb{E} [X_1^3] \right)^q \right]. \end{aligned}$$

Let  $X'_1, X'_2, \dots, X'_n$  be another set of  $n$  i.i.d. samples. We find

$$\begin{aligned} \mathbb{E} \left[ \left( \frac{1}{n} \sum_{i=1}^n X_i^3 - \mathbb{E} [X_1^3] \right)^q \right] &= \mathbb{E} \left[ \left( \frac{1}{n} \sum_{i=1}^n X_i^3 - \frac{1}{n} \sum_{i=1}^n \mathbb{E} [X_i^3] \right)^q \right] \\ &\stackrel{(a)}{\leq} \mathbb{E}_{X_i, X'_i} \left[ \left( \frac{1}{n} \sum_{i=1}^n (X_i^3 - X_i'^3) \right)^q \right] \\ &\stackrel{(b)}{\leq} \mathbb{E}_{X_i, X'_i, \sigma_i} \left[ \left( \frac{1}{n} \sum_{i=1}^n \sigma_i (X_i^3 - X_i'^3) \right)^q \right] \\ &\stackrel{(c)}{\leq} \mathbb{E}_{X_i, X'_i, \sigma_i} \left[ 2^{q-1} \left( \frac{1}{n} \sum_{i=1}^n \sigma_i X_i^3 \right)^q + 2^{q-1} \left( \frac{1}{n} \sum_{i=1}^n \sigma_i X_i'^3 \right)^q \right] \\ &= \left( \frac{2}{n} \right)^q \mathbb{E}_{X_i, \sigma_i} \left[ \left( \sum_{i=1}^n \sigma_i X_i^3 \right)^q \right], \end{aligned}$$

where (a) and (c) follow from Jensen's inequality. In step (b), we introduce Rademacher sequence  $\sigma_1, \sigma_2, \dots, \sigma_n$ , i.e.,  $\mathbb{P}(\sigma_i = 1) = \mathbb{P}(\sigma_i = -1) = 0.5$ . To ease notation, we let  $Z_i := \sigma_i X_i$ . So  $\sigma_i X_i^3 = Z_i^3$  and  $Z_i$  is still sub-Gaussian with norm  $K$ . It thus remains to bound  $\mathbb{E}[(\sum_{i=1}^n Z_i^3)^q]$ . Note that  $Z_i$  has symmetric distribution around 0, so  $\mathbb{E}[Z_i^a] = 0$  for any odd integer  $a$ . Accordingly, we have

$$\mathbb{E} \left[ \left( \sum_{i=1}^n Z_i^3 \right)^q \right] = \sum_{q_1 + \dots + q_n = q/2} \prod_{i=1}^n \mathbb{E} [Z_i^{6q_i}] \leq \sum_{q_1 + \dots + q_n = q/2} \prod_{i=1}^n (K \sqrt{6q_i})^{6q_i},$$

where the last inequality follows from the basic property that if  $X$  is sub-Gaussian random variable with norm  $K$ , then  $(\mathbb{E}[|X|^q])^{1/q} \leq K\sqrt{q}$  for all  $q > 1$ . Since all  $q_i \leq q/2$ , we have

$$\mathbb{E} \left[ \left( \sum_{i=1}^n Z_i^3 \right)^q \right] \leq \binom{q/2 + n - 1}{q/2} (K\sqrt{3q})^{3q} \leq \left( \frac{(q/2 + n - 1)e}{q/2} \right)^{q/2} (K\sqrt{3q})^{3q}.$$

Putting all pieces together, we have

$$\mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n X_i^3 - \mathbb{E}[X_1^3] \right| > t \right) \leq \left( \frac{18K^3 q \sqrt{q + 2n}}{nt} \right)^q.$$

Setting  $q = \lceil \log(1/\delta) \rceil$ ,  $t = 18eK^3 \frac{\sqrt{\log^3(1/\delta) + 2\log^2(1/\delta)n}}{n}$  completes the proof.  $\square$

**Lemma 12.** For any symmetric 3-way tensor  $\mathbf{T} \in \mathbb{R}^{p \times p \times p}$ ,

$$\sup_{\mathbf{u}, \mathbf{v}, \mathbf{w} \in \mathbb{S}^{p-1}} |\mathbf{T}(\mathbf{u}, \mathbf{v}, \mathbf{w})| \leq 9 \|\mathbf{T}\|_{op}.$$

*Proof.* For any  $\mathbf{u}, \mathbf{v}, \mathbf{w} \in \mathbb{S}^{p-1}$ , we have

$$\begin{aligned} 2|\mathbf{T}(\mathbf{u}, \mathbf{v}, \mathbf{w})| &= |\mathbf{T}(\mathbf{u}, \mathbf{v}, \mathbf{w}) + \mathbf{T}(\mathbf{v}, \mathbf{u}, \mathbf{w})| = |\mathbf{T}(\mathbf{u} + \mathbf{v}, \mathbf{u} + \mathbf{v}, \mathbf{w}) - \mathbf{T}(\mathbf{u}, \mathbf{u}, \mathbf{w}) - \mathbf{T}(\mathbf{v}, \mathbf{v}, \mathbf{w})| \\ &\leq |\mathbf{T}(\mathbf{u} + \mathbf{v}, \mathbf{u} + \mathbf{v}, \mathbf{w})| + |\mathbf{T}(\mathbf{u}, \mathbf{u}, \mathbf{w})| + |\mathbf{T}(\mathbf{v}, \mathbf{v}, \mathbf{w})| \leq 6 \sup_{\mathbf{a}, \mathbf{b} \in \mathbb{S}^{p-1}} |\mathbf{T}(\mathbf{a}, \mathbf{a}, \mathbf{b})|, \end{aligned}$$

where the first step holds because  $\mathbf{T}$  is symmetric. Moreover, for any  $\mathbf{u}, \mathbf{v} \in \mathbb{S}^{p-1}$ , we have

$$6\mathbf{T}(\mathbf{u}, \mathbf{u}, \mathbf{v}) = |\mathbf{T}(\mathbf{u} + \mathbf{v}, \mathbf{u} + \mathbf{v}, \mathbf{u} + \mathbf{v}) + \mathbf{T}(\mathbf{v} - \mathbf{u}, \mathbf{v} - \mathbf{u}, \mathbf{v} - \mathbf{u}) - 2\mathbf{T}(\mathbf{v}, \mathbf{v}, \mathbf{v})| \leq 18 \|\mathbf{T}\|_{op}.$$

Combining the above two inequalities leads to

$$\sup_{\mathbf{u}, \mathbf{v}, \mathbf{w} \in \mathbb{S}^{p-1}} |\mathbf{T}(\mathbf{u}, \mathbf{v}, \mathbf{w})| \leq 3 \sup_{\mathbf{u}, \mathbf{v} \in \mathbb{S}^{p-1}} |\mathbf{T}(\mathbf{u}, \mathbf{u}, \mathbf{v})| \leq 9 \|\mathbf{T}\|_{op}.$$

$\square$

**Lemma 13** (Conditional Mean Deviation). Let  $X \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ ,  $Z \sim \mathcal{N}(0, 1)$ , and assume  $X$  and  $Z$  are independent. For any  $\tau_1, \tau_2 \geq 1$ ,  $\mathbf{v} \in \mathbb{S}^{p-1}$ , we define event  $\mathcal{E} := \{|\langle X, \mathbf{v} \rangle| \leq \tau_1, |Z| \leq \tau_2\}$ . For any  $a, b > 0$ , let  $Y := a \cdot \langle X, \mathbf{v} \rangle + b \cdot Z$ . There exists constant  $C$  such that the following inequalities hold.

1.

$$\|\mathbb{E}[Y^3 X \mid \mathcal{E}] - \mathbb{E}[Y^3 X]\|_2 \leq C(a^3 + ab^2) \left( \tau_1^3 e^{-\tau_1^2/2} + \tau_1 \tau_2 e^{-\tau_1^2/2 - \tau_2^2/2} \right). \quad (62)$$

2.

$$\|\mathbb{E}[Y^2 X \otimes X \mid \mathcal{E}] - \mathbb{E}[Y^2 X \otimes X]\|_{op} \leq C(a^2 + b^2) \left( \tau_1^3 e^{-\tau_1^2/2} + \tau_1 \tau_2 e^{-\tau_1^2/2 - \tau_2^2/2} \right). \quad (63)$$

3.

$$\|\mathbb{E}[Y^3 X \otimes X \otimes X \mid \mathcal{E}] - \mathbb{E}[Y^3 X \otimes X \otimes X]\|_{op} \leq C(a^3 + ab^2) \left( \tau_1^5 e^{-\tau_1^2/2} + \tau_1^3 \tau_2 e^{-\tau_1^2/2 - \tau_2^2/2} \right). \quad (64)$$

*Proof.* 1. There exists  $\mathbf{u} \in \mathbb{S}^{p-1}$  such that

$$\delta_1 := \|\mathbb{E}[Y^3 X \mid \mathcal{E}] - \mathbb{E}[Y^3 X]\|_2 = \mathbb{E}[Y^3 \langle X, \mathbf{u} \rangle \mid \mathcal{E}] - \mathbb{E}[Y^3 \langle X, \mathbf{u} \rangle]$$

Due to the rotation invariance of spherical Gaussian vector, without loss of generality, we can simply assume  $\mathbf{v} = \mathbf{e}_1$  and  $\mathbf{u} = c\mathbf{e}_1 + d\mathbf{e}_2$ , where  $c^2 + d^2 = 1$ . Let  $X = (X_1, X_2, \dots, X_p)^\top$ . Using the symmetricity of  $X_1, Z, X_2$  when conditioning on  $\mathcal{E}^c$ , we have

$$\mathbb{E}[Y^3 \langle X, \mathbf{u} \rangle \mid \mathcal{E}] = \mathbb{E}[(aX_1 + bZ)^3(cX_1 + dX_2) \mid \mathcal{E}] = \mathbb{E}[a^3 c X_1^4 + 3ab^2 c X_1^2 Z^2 \mid \mathcal{E}] \lesssim a^3 |c| + ab^2 |c|.$$

Note that  $X_1, Z, X_2$  are also symmetric when conditioning on  $\mathcal{E}^c$ , we thus obtain

$$\begin{aligned} \mathbb{E}[Y^3 \langle X, \mathbf{u} \rangle \mid \mathcal{E}^c] \mathbb{P}(\mathcal{E}^c) &= \mathbb{E}[a^3 c X_1^4 + 3ab^2 c X_1^2 Z^2 \mid \mathcal{E}^c] \mathbb{P}(\mathcal{E}^c) \\ &\lesssim a^3 |c| \tau_1^3 e^{-\tau_1^2/2} + ab^2 |c| \tau_1 \tau_2 e^{-\tau_1^2/2 - \tau_2^2/2}, \end{aligned}$$

where the last inequality follows from Lemma 14. Now we turn to  $\delta_1$ . We find

$$\begin{aligned} \delta_1 &= \mathbb{E}[Y^3 \langle X, \mathbf{u} \rangle \mid \mathcal{E}] - \mathbb{E}[Y^3 \langle X, \mathbf{u} \rangle \mid \mathcal{E}] \mathbb{P}(\mathcal{E}) - \mathbb{E}[Y^3 \langle X, \mathbf{u} \rangle \mid \mathcal{E}^c] \mathbb{P}(\mathcal{E}^c) \\ &\leq |\mathbb{E}[Y^3 \langle X, \mathbf{u} \rangle \mid \mathcal{E}]| \mathbb{P}(\mathcal{E}^c) + |\mathbb{E}[Y^3 \langle X, \mathbf{u} \rangle \mid \mathcal{E}^c]| \mathbb{P}(\mathcal{E}^c) \\ &\lesssim (a^3 |c| + 3ab^2 |c|) e^{-\tau_1^2/2 - \tau_2^2/2} + a^3 |c| \tau_1^3 e^{-\tau_1^2/2} + ab^2 |c| \tau_1 \tau_2 e^{-\tau_1^2/2 - \tau_2^2/2} \\ &\lesssim (a^3 + ab^2) \left( \tau_1^3 e^{-\tau_1^2/2} + \tau_1 \tau_2 e^{-\tau_1^2/2 - \tau_2^2/2} \right). \end{aligned}$$

2. There exists  $\mathbf{u} \in \mathbb{S}^{p-1}$  such that

$$\delta_2 := \|\mathbb{E}[Y^2 X \otimes X \mid \mathcal{E}] - \mathbb{E}[Y^2 X \otimes X]\|_{op} = \mathbb{E}[Y^2 \langle X, \mathbf{u} \rangle^2 \mid \mathcal{E}] - \mathbb{E}[Y^2 \langle X, \mathbf{u} \rangle^2].$$

Using the same simplification argument in (a), we have

$$\begin{aligned} \mathbb{E}[Y^2 \langle X, \mathbf{u} \rangle^2 \mid \mathcal{E}] &= \mathbb{E}[(aX_1 + bZ)^2(cX_1 + dX_2)^2 \mid \mathcal{E}] \\ &= \mathbb{E}[a^2 c^2 X_1^4 + b^2 c^2 X_1^2 Z^2 + a^2 d^2 X_1^2 X_2^2 + b^2 d^2 X_2^2 Z^2 \mid \mathcal{E}] \lesssim a^2 + b^2. \end{aligned}$$

Applying Lemma 14 again leads to

$$\begin{aligned} \mathbb{E}[Y^2 \langle X, \mathbf{u} \rangle^2 \mid \mathcal{E}^c] \mathbb{P}(\mathcal{E}^c) &\lesssim a^2 c^2 \tau_1^3 e^{-\tau_1^2/2} + b^2 c^2 \tau_1 \tau_2 e^{-\tau_1^2/2 - \tau_2^2/2} + a^2 d^2 \tau_1 e^{-\tau_1^2/2} + b^2 d^2 \tau_2 e^{-\tau_2^2/2} \\ &\lesssim (a^2 + b^2) \left( \tau_1^3 e^{-\tau_1^2/2} + \tau_1 \tau_2 e^{-\tau_1^2/2 - \tau_2^2/2} \right). \end{aligned}$$

Overall, we have

$$\begin{aligned}\delta_2 &= \mathbb{E} [Y^2 \langle X, \mathbf{u} \rangle^2 \mid \mathcal{E}] - \mathbb{E} [Y^2 \langle X, \mathbf{u} \rangle^2 \mid \mathcal{E}] \mathbb{P}(\mathcal{E}) - \mathbb{E} [Y^2 \langle X, \mathbf{u} \rangle^2 \mid \mathcal{E}^c] \mathbb{P}(\mathcal{E}^c) \\ &\leq |\mathbb{E} [Y^2 \langle X, \mathbf{u} \rangle^2 \mid \mathcal{E}] \mathbb{P}(\mathcal{E}^c) + |\mathbb{E} [Y^2 \langle X, \mathbf{u} \rangle^2 \mid \mathcal{E}^c] \mathbb{P}(\mathcal{E})|. \\ &\lesssim (a^2 + b^2) \left( \tau_1^3 e^{-\tau_1^2/2} + \tau_1 \tau_2 e^{-\tau_1^2/2 - \tau_2^2/2} \right).\end{aligned}$$

3. There exists  $\mathbf{u} \in \mathbb{S}^{p-1}$  such that

$$\delta_3 := \|\mathbb{E} [Y^3 X \otimes X \otimes X \mid \mathcal{E}] - \mathbb{E} [Y^3 X \otimes X \otimes X]\|_{op} = \mathbb{E} [Y^3 \langle X, \mathbf{u} \rangle^3 \mid \mathcal{E}] - \mathbb{E} [Y^3 \langle X, \mathbf{u} \rangle^3].$$

Using the same simplification argument in (a), we have

$$\begin{aligned}\mathbb{E} [Y^3 \langle X, \mathbf{u} \rangle^3 \mid \mathcal{E}] &= \mathbb{E} [(aX_1 + bZ)^3 (cX_1 + dX_2)^3 \mid \mathcal{E}] \\ &= \mathbb{E} [a^3 c^3 X_1^6 + 3ab^2 c^3 X_1^4 Z^2 + 3a^3 cd^2 X_1^4 X_2^2 + 9ab^2 cd^2 X_1^2 X_2^2 Z^2 \mid \mathcal{E}] \lesssim a^3 + ab^2.\end{aligned}$$

Applying Lemma 14 again leads to

$$\begin{aligned}\mathbb{E} [Y^3 \langle X, \mathbf{u} \rangle^3 \mid \mathcal{E}^c] \mathbb{P}(\mathcal{E}^c) &\lesssim a^3 c^3 \tau_1^5 e^{-\tau_1^2/2} + ab^2 c^2 \tau_1^3 e^{-\tau_1^2/2} + a^3 cd^2 \tau_1^3 \tau_2 e^{-\tau_1^2/2 - \tau_2^2/2} + ab^2 cd^2 \tau_1 \tau_2 e^{-\tau_1^2/2 - \tau_2^2/2} \\ &\lesssim (a^3 + ab^2) \left( \tau_1^5 e^{-\tau_1^2/2} + \tau_1^3 \tau_2 e^{-\tau_1^2/2 - \tau_2^2/2} \right).\end{aligned}$$

Finally, we have

$$\begin{aligned}\delta_3 &= \mathbb{E} [Y^3 \langle X, \mathbf{u} \rangle^3 \mid \mathcal{E}] - \mathbb{E} [Y^3 \langle X, \mathbf{u} \rangle^3 \mid \mathcal{E}] \mathbb{P}(\mathcal{E}) - \mathbb{E} [Y^3 \langle X, \mathbf{u} \rangle^3 \mid \mathcal{E}^c] \mathbb{P}(\mathcal{E}^c) \\ &\leq |\mathbb{E} [Y^3 \langle X, \mathbf{u} \rangle^3 \mid \mathcal{E}] \mathbb{P}(\mathcal{E}^c) + |\mathbb{E} [Y^3 \langle X, \mathbf{u} \rangle^3 \mid \mathcal{E}^c] \mathbb{P}(\mathcal{E})|. \\ &\lesssim (a^3 + ab^2) \left( \tau_1^5 e^{-\tau_1^2/2} + \tau_1^3 \tau_2 e^{-\tau_1^2/2 - \tau_2^2/2} \right).\end{aligned}$$

□

**Lemma 14** (Conditional Moments of Gaussian). *Suppose  $X \sim \mathcal{N}(0, 1)$ . For any  $\tau > 0$  and positive integer  $a$ , we define*

$$m_a(\tau) := \mathbb{E} [X^a \mid |X| > \tau] \mathbb{P}(|X| > \tau).$$

*Then we have that for all  $a = 2, 4, 6, \dots$ , we have*

$$m_a(\tau) = (a - 1)m_{a-2}(\tau) + \sqrt{\frac{2}{\pi}} \tau^{a-1} e^{-\frac{\tau^2}{2}}.$$

*Proof.* The result follows from elementary calculation on Gaussian's probability density function. We omit the details. □

**Lemma 15** (Sub-Gaussianity). *Let  $X \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ . For any  $k$  fixed vectors  $\mathbf{u}_1, \dots, \mathbf{u}_k \in \mathbb{R}^p$ , we define event*

$$\mathcal{E} := \{|\langle X, \mathbf{u}_1 \rangle| \leq |\langle X, \mathbf{u}_j \rangle|, \text{ for all } j \in [k]\}.$$

(a) *Suppose  $\mathbb{P}(\mathcal{E}) \geq \tau > 0$ . There exists constant  $C$  that only depends on  $\tau$  such that for any fixed  $\mathbf{x} \in \mathbb{S}^{p-1}$ , we have that*

$$\mathbb{P}(|\langle X, \mathbf{x} \rangle| > t \mid \mathcal{E}) \leq e^{1-Ct^2}, \text{ for all } t > 0.$$

(b) *In general there exists constant  $C'$  such that for any fixed  $\mathbf{x} \in \mathbb{S}^{p-1}$ ,*

$$\mathbb{P}(|\langle X, \mathbf{x} \rangle| > t \mid \mathcal{E}) \leq e^{1 - \frac{C'}{4k \log(k+1)} t^2}, \text{ for all } t > 0.$$



*Proof.* Since  $X$  is sub-Gaussian random vector, equivalently there exists constant  $C$  such that for any fixed  $\mathbf{x} \in \mathbb{S}^{p-1}$ ,

$$\mathbb{P}\{|\langle X, \mathbf{x} \rangle| \geq t\} \leq 1 \wedge e^{1-Ct^2}, \text{ for all } t > 0.$$

(a) Note that

$$\begin{aligned} \mathbb{P}(|\langle X, \mathbf{x} \rangle| \geq t) &= \mathbb{P}(|\langle X, \mathbf{x} \rangle| \geq t \mid \mathcal{E}) \mathbb{P}(\mathcal{E}) + \mathbb{P}(|\langle X, \mathbf{x} \rangle| \geq t \mid \mathcal{E}^c) \mathbb{P}(\mathcal{E}^c) \\ &\geq \tau \cdot \mathbb{P}(|\langle X, \mathbf{x} \rangle| \geq t \mid \mathcal{E}). \end{aligned}$$

Hence,

$$\mathbb{P}(|\langle X, \mathbf{x} \rangle| \geq t \mid \mathcal{E}) \leq 1 \wedge \tau^{-1} e^{1-Ct^2} \leq 1 \wedge e^{1-C'(\tau)t^2},$$

where the last inequality holds for  $C'(\tau) = C(1 - \log \tau)^{-1}$ .

(b) Without loss of generality, we assume that  $\mathbf{u}_1, \dots, \mathbf{u}_k$  live in the subspace spanned by  $\mathbf{e}_1, \dots, \mathbf{e}_k$ . For any vector  $\mathbf{u} \in \mathbb{R}^p$ , we let  $\mathbf{u}_{[k]}$  be its sub-vector that contains the first  $k$  coordinates, and  $\mathbf{u}_\perp$  be its sub-vector that contains the rest coordinates. For any  $\mathbf{x} \in \mathbb{S}^{p-1}$ , we have

$$\begin{aligned} \mathbb{P}(|\langle X, \mathbf{x} \rangle| > t \mid \mathcal{E}) &\leq \mathbb{P}(|\langle X_{[k]}, \mathbf{x}_{[k]} \rangle| > t/2 \mid \mathcal{E}) + \mathbb{P}(|\langle X_\perp, \mathbf{x}_\perp \rangle| > t/2) \\ &\leq \mathbb{P}(|\langle X_{[k]}, \mathbf{x}_{[k]} \rangle| > t/2 \mid \mathcal{E}) + e^{1-Ct^2/4} \\ &\leq \mathbb{P}(\|X_{[k]}\|_2 > t/2 \mid \mathcal{E}) + e^{1-Ct^2/4}. \end{aligned} \tag{65}$$

Note that conditioning  $\mathcal{E}$  does not change the distribution of  $\|X_{[k]}\|_2$ . We thus have

$$\mathbb{P}(\|X_{[k]}\|_2 > t/2 \mid \mathcal{E}) = \mathbb{P}(\|X_{[k]}\|_2 > t/2) \leq \sum_{i \in [k]} \mathbb{P}\left(|X_i| \geq \frac{t}{2\sqrt{k}}\right) \leq k \cdot e^{1-Ct^2/(4k)}.$$

Combining (65) with the above inequality yields that

$$\mathbb{P}(|\langle X, \mathbf{x} \rangle| > t \mid \mathcal{E}) \leq 1 \wedge (k+1)e^{1-Ct^2/(4k)} \leq 1 \wedge e^{1-C_2(k)t^2},$$

where the last inequality holds by setting  $C_2(k) = \frac{C}{4k \log(k+1)}$ . □

## References

- [1] Animashree Anandkumar, Daniel Hsu, and Sham M Kakade. “A method of moments for mixture models and hidden Markov models”. In *arXiv preprint arXiv:1203.0683* (2012).
- [2] Animashree Anandkumar, Rong Ge, Daniel Hsu, Sham M. Kakade, and Matus Telgarsky. “Tensor decompositions for learning latent variable models”. In *The Journal of Machine Learning Research* 15.1 (2014), pp. 2773–2832.
- [3] Animashree Anandkumar, Sham M Kakade, Dean P Foster, Yi-Kai Liu, and Daniel Hsu. *Two svds suffice: Spectral decompositions for probabilistic topic modeling and latent dirichlet allocation*. Tech. rep. 2012.
- [4] Sivaraman Balakrishnan, Martin J. Wainwright, and Bin Yu. “Statistical guarantees for the EM algorithm: From population to sample-based analysis”. In *arXiv preprint arXiv:1408.2156* (2014).

- [5] Emmanuel J. Candès, Xiaodong Li, and Mahdi Soltanolkotabi. “Phase retrieval via Wirtinger flow: Theory and algorithms”. In *IEEE Transactions on Information Theory* 61.4 (2015), pp. 1985–2007.
- [6] Arun Chaganty and Percy Liang. “Spectral Experts for Estimating Mixtures of Linear Regressions”. In *International Conference on Machine Learning (ICML)*. 2013.
- [7] Yudong Chen and Martin J Wainwright. “Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees”. In *arXiv preprint arXiv:1509.03025* (2015).
- [8] Yudong Chen, Xinyang Yi, and Constantine Caramanis. “A Convex Formulation for Mixed Regression with Two Components: Minimax Optimal Rates.” In *COLT*. 2014, pp. 560–604.
- [9] Yuxin Chen and Emmanuel J. Candès. “Solving random quadratic systems of equations is nearly as easy as solving linear systems”. In *Advances in Neural Information Processing Systems*. 2015, pp. 739–747.
- [10] Partha Deb and Ann M. Holmes. “Estimates of use and costs of behavioural health care: a comparison of standard and finite mixture models”. In *Health Economics* 9.6 (2000), pp. 475–489.
- [11] Bettina Grün and Friedrich Leisch. “Applications of finite mixtures of regression models”. In *URL: <http://cran.r-project.org/web/packages/flexmix/vignettes/regression-examples.pdf>* (2007).
- [12] Daniel Hsu and Sham M. Kakade. “Learning Gaussian Mixture Models: Moment Methods and Spectral Decompositions”. In *CoRR* abs/1206.5766 (2012).
- [13] Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. “Adaptive mixtures of local experts”. In *Neural computation* 3.1 (1991), pp. 79–87.
- [14] Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi. “Low-rank matrix completion using alternating minimization”. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*. ACM. 2013, pp. 665–674.
- [15] Karl Pearson. “Contributions to the mathematical theory of evolution”. In *Philosophical Transactions of the Royal Society of London. A* (1894), pp. 71–110.
- [16] Hanie Sedghi and Anima Anandkumar. “Provable Tensor Methods for Learning Mixtures of Generalized Linear Models”. In *arXiv preprint arXiv:1412.3046* (2014).
- [17] Hanie Sedghi, Majid Janzamin, and Anima Anandkumar. “Provable Tensor Methods for Learning Mixtures of Generalized Linear Models”. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*. 2016, pp. 1223–1231.
- [18] Ruoyu Sun and Zhi-Quan Luo. “Guaranteed matrix completion via nonconvex factorization”. In *Foundations of Computer Science (FOCS), 2015 IEEE 56th Annual Symposium on*. IEEE. 2015, pp. 270–289.
- [19] Yuekai Sun, Stratis Ioannidis, and Andrea Montanari. “Learning Mixtures of Linear Classifiers.” In *ICML*. 2014, pp. 721–729.
- [20] Aad W Van Der Vaart and Jon A Wellner. *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer, 1996.

- [21] Roman Vershynin. “Introduction to the non-asymptotic analysis of random matrices”. In *Arxiv preprint arxiv:1011.3027* (2010).
- [22] Kert Viece and Barbara Tong. “Modeling with Mixtures of Linear Regressions”. In *Statistics and Computing* 12.4 (2002). ISSN: 0960-3174. URL: <http://dx.doi.org/10.1023/A%3A1020779827503>.
- [23] Michel Wedel and Wayne S DeSarbo. “A mixture likelihood approach for generalized linear models”. In *Journal of Classification* 12.1 (1995), pp. 21–55.
- [24] Michel Wedel and Wagner A Kamakura. *Market segmentation: Conceptual and methodological foundations*. Vol. 8. Springer Science & Business Media, 2012.
- [25] Xinyang Yi, Constantine Caramanis, and Sujay Sanghavi. “Alternating Minimization for Mixed Linear Regression.” In *ICML*. 2014, pp. 613–621.
- [26] Yuchen Zhang, Xi Chen, Denny Zhou, and Michael I Jordan. “Spectral methods meet EM: A provably optimal algorithm for crowdsourcing”. In *Advances in neural information processing systems*. 2014, pp. 1260–1268.