

# Overcomplete Order-3 Tensor Decomposition, Blind Deconvolution, and Gaussian Mixture Models\*

Haolin Chen<sup>†</sup> and Luis Rademacher<sup>†</sup>

**Abstract.** We propose a new algorithm for tensor decomposition, based on the simultaneous diagonalization algorithm, and apply our new algorithmic ideas to blind deconvolution and Gaussian mixture models. Our first contribution is a simple and efficient algorithm to decompose certain symmetric overcomplete order-3 tensors, that is, three dimensional arrays of the form  $T = \sum_{i=1}^n a_i \otimes a_i \otimes a_i$  where the  $a_i$ s are not linearly independent. Our algorithm comes with a detailed robustness analysis. Our second contribution builds on top of our tensor decomposition algorithm to expand the family of Gaussian mixture models whose parameters can be estimated efficiently. These ideas are also presented in a more general framework of blind deconvolution that makes them applicable to mixture models of identical but very general distributions, including all centrally symmetric distributions with finite 6th moment.

**Key words.** tensor decomposition, blind deconvolution, Gaussian mixture models

**AMS subject classifications.** 15A69, 62H30, 68T09, 68W20

**DOI.** 10.1137/21M1399415

**1. Introduction.** Tensor decomposition is a basic tool in data analysis. The *order-3 (symmetric) tensor decomposition problem*<sup>1</sup> can be stated as follows: Given an order-3 tensor  $T = \sum_{i=1}^n a_i \otimes a_i \otimes a_i$ , recover the vectors  $a_i \in \mathbb{R}^d$ . The problem is *undercomplete* if the  $a_i$ s are linearly independent; otherwise, it is *overcomplete*. Two problems in data analysis motivate us here to study tensor decomposition: blind deconvolution and Gaussian mixture models (GMMs).

A *deconvolution* problem can be formulated as follows: We have a  $d$ -dimensional random vector

$$(1.1) \quad X = Z + \eta,$$

where  $Z$  and  $\eta$  are independent random vectors. Given samples from  $X$ , the goal is to determine the distribution of  $Z$ . We call it *blind deconvolution* when the distribution of  $\eta$  is unknown; otherwise, it is *nonblind*. It is called *deconvolution* because the distribution of  $X$  is the convolution of the distributions of  $Z$  and  $\eta$ .

\*Received by the editors February 18, 2021; accepted for publication (in revised form) November 29, 2021; published electronically March 9, 2022.

<https://doi.org/10.1137/21M1399415>

**Funding:** This work was supported by National Science Foundation grants CCF-1657939, CCF-1422830, CCF-2006994, and CCF-1934568.

<sup>†</sup>Department of Mathematics, University of California, Davis, Davis, CA 95616 USA ([hlnchen@ucdavis.edu](mailto:hlnchen@ucdavis.edu), [lrademac@ucdavis.edu](mailto:lrademac@ucdavis.edu)).

<sup>1</sup>“Tensor decomposition” here is a shorthand for a specific kind of tensor decomposition, the symmetric tensor rank decomposition of a symmetric tensor. See [section 2](#) for a discussion.

The following *mixture model parameter estimation problem* can be recast as a blind deconvolution problem: Let  $X$  be a  $d$ -dimensional random vector distributed as the following mixture model: First sample  $i$  from  $[d]$ , each value with probability  $w_i$  ( $w_i > 0, \sum_i w_i = 1$ ), and then let  $X = \mu_i + \eta$ , where  $\eta$  is a given  $d$ -dimensional random vector and  $\mu_i \in \mathbb{R}^d$ . The estimation problem is to estimate  $\mu_i$ s and  $w_i$ s from samples of  $X$ . It is a deconvolution problem  $X = Z + \eta$  when  $Z$  follows the discrete distribution equal to  $\mu_i$  with probability  $w_i$ . It is blind when the distribution of  $\eta$  is unknown.

The *GMM parameter estimation problem* can be described as follows: Let  $X \in \mathbb{R}^d$  be a random vector with density function  $x \mapsto \sum_{i=1}^k w_i f_i(x)$ , where  $w_i > 0, \sum_i w_i = 1$ , and  $f_i$  is the Gaussian density function with mean  $\mu_i \in \mathbb{R}^d$  and covariance matrix  $\Sigma_i \in \mathbb{R}^{d \times d}$ . GMM parameter estimation is the following algorithmic question: Given i.i.d. samples from  $X$ , estimate  $w_i$ s,  $\mu_i$ s, and  $\Sigma_i$ s.

The GMM parameter estimation problem is a deconvolution problem when the covariance matrices of the components are the same, namely  $\Sigma_i = \Sigma$ . Specifically,  $X = Z + \eta$ , where  $Z$  follows a discrete distribution taking value  $\mu_i$  with probability  $w_i, i = 1, \dots, k$ , and  $\eta$  is Gaussian with mean 0 and covariance  $\Sigma$ . It is blind if  $\Sigma$  is unknown.

While the undercomplete tensor decomposition problem is well understood (based on algorithmic techniques such as the tensor power method and the simultaneous diagonalization algorithm [28]), the overcomplete regime is much more challenging [21, Chapter 7]. Within the overcomplete case, there are fewer techniques available for the order-3 case than there are for higher order ones [21, section 7.3]. We discuss some of these techniques and challenges below (subsection 1.2).

**1.1. Our results.**

*Overcomplete tensor decomposition.* We propose an algorithm based on the simultaneous diagonalization algorithm for overcomplete tensor decomposition. Our informal claim is as follows.

**Claim 1.1 (informal statement of Theorem 3.1).** *Given a symmetric order-3 tensor  $T = \sum_{i=1}^{d+k} a_i^{\otimes 3} \in \mathbb{R}^{d \times d \times d}$  and when any  $d$ -subset of the  $a_i$ s is linearly independent, there is a randomized algorithm that recovers  $a_i$ s within  $\varepsilon$  error and with expected running time polynomial in  $d^k, 1/\varepsilon^k$ , and natural conditioning parameters.*

Note that our goal is to show that the running time has polynomial dependence in that sense and the error has inverse polynomial dependence, but we do not optimize the degrees of the polynomials. Even though the algorithm is exponential in  $k$ , the case  $k = 1$  already makes possible a new GMM result (see below).

Our proposed algorithm (Algorithm 2) and its analysis (Theorem 3.1) are stronger than Claim 1.1 in two important ways: It is robust in the sense that it approximates the  $a_i$ s even when the input is a tensor that is  $\varepsilon'$ -close to  $T$ . Also, it turns out that parameter  $k$  above, the number of  $a_i$ s beyond the dimension  $d$ , is not the best notion of overcompleteness. In our result, the tensor is of the form  $T = \sum_{i=1}^{r+k} a_i \otimes a_i \otimes a_i$ , where  $r$  is the robust Kruskal rank of  $a_i$ s (informally the maximum  $r$  such that any  $r$ -subset is well-conditioned; see Definition 2.1), so that  $k$  is the number of components above the robust Kruskal rank. Thus, our analysis also applies when the Kruskal rank is less than  $d$ .

We show the relevance of our algorithm with new results in two applications: blind deconvolution and learning GMMs. Our approach to these applications is decomposing the 3rd cumulant tensor of the underlying distribution, yielding efficient algorithms with [Algorithm 2](#). We restrict our study to the 3rd cumulant tensor (instead of higher order tensors that could be used in those applications) mainly because (1) the decomposition of the 3rd order tensor is less understood compared to tensors of higher order; (2) the cost of estimating higher order tensors is in principle higher; (3) the decomposition of the 3rd cumulant tensor yields easily an algorithm for those applications.

**Blind deconvolution.** We provide an efficient algorithm for the following blind deconvolution problem.

**Claim 1.2 (informal statement of [Theorem 5.3](#)).** *Let  $X = Z + \eta$  be a random vector as in (1.1), where  $Z$  is a  $d$ -dimensional discrete distribution supported on  $d$  points and  $\eta$  has zero mean, zero 3rd moment, and finite 6th moment. Suppose  $Z$  satisfies a natural nondegeneracy condition ([Assumption 5.1](#)). Then there is a randomized algorithm that, with probability  $1 - \delta$  over the randomness in the samples, recovers  $Z$  within  $\varepsilon$  error. The expected running time and sample complexity are polynomial in  $d, \varepsilon^{-1}, \delta^{-1}$ , and natural condition parameters.*

Equivalently, it can solve the mixture model parameter estimation problem above under the same conditions ([Algorithm 3](#) and [Theorem 5.3](#)).

**GMM.** We show an efficient algorithm for the following GMM parameter estimation problem.

**Claim 1.3 (informal statement of [Theorem 6.1](#)).** *Given samples from a  $d$ -dimensional mixture of  $d$  identical and not necessarily spherical Gaussians with unknown parameters  $w_i, \mu_i, \Sigma$  satisfying a natural nondegeneracy condition ([Assumption 5.1](#)), there is a randomized algorithm that with probability  $1 - \delta$  over the randomness in the samples estimates all parameters within  $\varepsilon$  error. The expected running time and sample complexity are polynomial in  $d, \varepsilon^{-1}, \delta^{-1}$ , and natural conditioning parameters.*

It may seem as if the last two contributions (blind deconvolution and GMM) could be attacked with standard *undercomplete* tensor decomposition techniques, given that the number of components is equal to the ambient dimension and therefore they could be linearly independent. It is not clear how that could actually happen, as the nonspherical unknown covariance seems to make standard approaches inapplicable and our contribution is a formulation that involves an overcomplete tensor decomposition and uses our overcomplete tensor decomposition algorithm in an essential way.

**Organization of the paper.** In [section 2](#), we introduce the notation and preliminaries needed in the paper. In [section 3](#), we present the proposed algorithm ([Algorithm 2](#)) for overcomplete tensor decomposition and its analysis ([Theorem 3.1](#)), as well as our high level proof ideas. We implement our proof ideas for [Theorem 3.1](#) in [section 4](#). In [sections 5](#) and [6](#), we provide the algorithms ([Algorithms 3](#) and [4](#)) and their analyses ([Theorems 5.3](#) and [6.1](#)) for the two applications mentioned above.

**1.2. Related work.** Among basic tensor decomposition techniques for the undercomplete case we have *tensor power iteration* (see [\[21\]](#), for example) and *the simultaneous diagonalization algorithm* [\[28\]](#). Tensor power iteration is more robust than the simultaneous diagonal-

ization algorithm, while the simultaneous diagonalization algorithm can be applied more generally: Tensor power iteration is mainly an algorithm for orthogonal tensors (orthogonal  $a_i$ s) and the general case with additional information, while the simultaneous diagonalization algorithm can decompose the general case without additional information. Our contributions below are based on the simultaneous diagonalization algorithm because of this additional power. The robustness of the simultaneous diagonalization algorithm is studied in several papers; our analysis builds on top of [5, 16].

For the overcomplete regime, we have algorithms such as FOABI [27] and the works [1, 2, 5, 14, 15, 18, 30].

Many techniques for the overcomplete case only make sense for orders 4 and higher or have weaker guarantees in the order-3 case. For example, some techniques use the fact that a  $d \times d \times d \times d$  tensor can be seen as a  $d^2 \times d^2$  matrix (and similarly for orders higher than 4), while no equally useful operation is available for order-3 tensors. Nevertheless, there are several results about decomposition in the order-3 case that are relevant to our work.

Kruskal gave a sufficient condition for unique decomposition [25]. A robust version of Kruskal's uniqueness and an algorithm running in time exponential in the number of components is given by Bhaskara, Charikar, and Vijayaraghavan [6]. Kruskal's uniqueness and its robust version will be building blocks of our results.

Anandkumar, Ge, and Janzamin [1, 2] develop an algorithm based on tensor power iteration that holds up to  $n \leq \beta d$  components for any  $\beta \geq 1$ . The running time is polynomial in  $d$  and exponential in  $\beta$ . The analysis therein requires assumptions such as incoherent components ( $\max_{i \neq j} |\langle a_i, a_j \rangle|$  is upper bounded inverse polynomially) to guarantee the convergence of tensor power iteration.

Recently, sum-of-squares based algorithms are proposed to solve *random* overcomplete order-3 tensor decomposition, where all components are Gaussian or uniform on the sphere. The algorithm in Ge and Ma [14] can decompose up to  $n \leq d^{3/2}/(\log d)^{O(1)}$  components and runs in  $n^{O(\log n)}$  time. Furthermore, the algorithm in Hopkins et al. [19] has the running time boosted to  $O(nd^{1+\omega})$ , where  $\omega < 2.4$  is the exponent of matrix multiplication. At the same time, the algorithm suffers a sacrifice of handling only up to  $d^{4/3}/(\log d)^{O(1)}$  components. Ma, Shi, and Steurer [31] further improved the number of components to  $d^{3/2}/(\log d)^{O(1)}$  and the running time to polynomial in  $n$ . These works provide an understanding of overcomplete *random* tensor decomposition in the regime of high rank, while ours focuses on the mildly overcomplete case and our main assumption follows naturally from the weaker assumptions in Kruskal's theorem.

Among works closest to ours, [11, 12] propose an algorithm that is efficient in the mildly overcomplete case for overcomplete order-3 tensor decomposition under natural nondegeneracy conditions. Though our results have similar assumptions and computational cost compared to [11, 12], our algorithm is comparatively a very simple randomized algorithm and we provide a rigorous robustness analysis.

Blind deconvolution-type problems have a long history in signal processing and specifically in image processing as a deblurring technique (see, e.g., [29]). The idea of using higher order moments in blind identification problems is standard too in signal processing, specifically in independent component analysis (see, e.g., [7, 8]). Our model (1.1) is somewhat different but very natural and inspired by mixture models.

With respect to GMMs, we are interested in parameter estimation in high dimension with no separation assumption (i.e., the means  $\mu_i$  can be arbitrarily close). Among the most relevant results in this context we have the following polynomial time algorithms: [20] for linearly independent means and spherical components (each  $\Sigma_i$  is a multiple of the identity); [3] for  $O(d^c)$  components with identical and known covariance  $\Sigma$ ; [5] for  $O(d^c)$  components with each  $\Sigma_i$  being diagonal in the smoothed analysis sense; [17, section 7], [16] for linearly independent means and spherical components in the presence of Gaussian noise; and [13] for a general GMM with  $O(\sqrt{d})$  components in the sense of smoothed analysis. Our algorithm expands the family of GMMs for which efficient parameter estimation is possible. It does not require prior knowledge of the covariance matrix, unlike [3], and can handle more components ( $d$  components) than [13] at the price of assuming that all covariance matrices are identical. With respect to recent results on clustering-based algorithms [10, 22], we consider these works incomparable to ours since clustering-based algorithms typically require some separation assumptions in the parameters.

**2. Notation and preliminaries.** For clarity of exposition, we analyze our algorithms in a computational model where we assume arithmetic operations between real numbers take constant time. We use the notation  $\text{poly}(\cdot)$  to denote a fixed polynomial that is nondecreasing in every argument. See [17, section 5.3] for a discussion of the complexity of the simultaneous diagonalization algorithm.

For  $n \in \mathbb{N}$ , let  $[n] = \{1, \dots, n\}$ . The unit sphere in  $\mathbb{R}^d$  is denoted by  $\mathcal{S}^{d-1}$ .

**Matrices and vectors.** For a matrix  $A \in \mathbb{R}^{m \times n}$ , we denote by  $\sigma_i(A)$  its  $i$ th largest singular value, by  $A^\dagger$  its Moore–Penrose pseudoinverse, and by  $\kappa(A) = \sigma_1(A)/\sigma_{\min(m,n)}(A)$  its condition number. Let  $\text{vec}(A) \in \mathbb{R}^{mn}$  denote the vector obtained by stacking all columns of  $A$ . Denote by  $\text{diag}(a)$  the diagonal matrix with diagonal entries from  $a$ , where  $a$  is a (column) vector. Let  $\|\cdot\|_2$  denote the spectral norm of a matrix and  $\|\cdot\|_F$  the Frobenius norm of a matrix.

In  $\mathbb{R}^d$ , we denote by  $\langle a, b \rangle$  the inner product of two vectors  $a, b$ . Let  $\hat{a} = a/\|a\|_2$ . For a set of vectors  $\{a_1, a_2, \dots, a_n\}$ , we denote their linear span by  $\text{span}\{a_1, \dots, a_n\}$ . We use  $[a_1, a_2, \dots, a_n]$  to denote the matrix containing  $a_i$ s as columns. If  $A = [a_1, a_2, \dots, a_n]$ , we have  $\hat{A} = [\hat{a}_1, \hat{a}_2, \dots, \hat{a}_n]$  and  $\tilde{A}$  follows a similar definition. We denote by  $A_m \in \mathbb{R}^{d \times m}$  the matrix  $[a_1, a_2, \dots, a_m]$  for some  $m < n$  and by  $A_{>m} \in \mathbb{R}^{d \times (n-m)}$  the matrix  $[a_{m+1}, \dots, a_n]$ . We say the matrix  $A$  is  $\rho$ -bounded if  $\max_{i \in [n]} \|a_i\|_2 \leq \rho$ . Given a vector  $a \in \mathbb{R}^d$  or a diagonal matrix  $D \in \mathbb{R}^{d \times d}$ , for  $r \in \mathbb{R}$ , notations  $a^r$  and  $D^r$  are used for entrywise power.

**Definition 2.1** (see [25, 6]). *Let  $A \in \mathbb{R}^{m \times n}$ . The Kruskal rank of  $A$ , denoted by  $\text{K-rank}(A)$ , is the maximum  $k \in [n]$  such that any  $k$  columns of  $A$  are linearly independent. Let  $\tau > 0$ . The robust Kruskal rank (with threshold  $\tau$ ) of  $A$ , denoted by  $\text{K-rank}_\tau(A)$ , is the maximum  $k \in [n]$  such that for any subset  $S \subseteq [n]$  of size  $k$  we have  $\sigma_k(A_S) \geq 1/\tau$ .*

**Tensors.** A symmetric order-3 tensor is a three dimensional array  $T \in \mathbb{R}^{d \times d \times d}$  such that entry  $T_{ijk}$  is invariant under permutation of indices  $i, j, k$ . For a symmetric order-3 tensor  $T \in \mathbb{R}^{d \times d \times d}$  and a vector  $x \in \mathbb{R}^d$ , let  $T_x$  denote the matrix  $\sum_{i,j,k=1}^d T_{ijk} x_i e_j e_k^\top \in \mathbb{R}^{d \times d}$ . Let  $a^{\otimes 3}$  be a shorthand for  $a \otimes a \otimes a$ . For this paper, the *rank* of a symmetric tensor  $T$  is a shorthand for its *symmetric rank*, namely the minimal  $n$  such that  $T = \sum_{i=1}^n a_i^{\otimes 3}$ . For a rank  $n$  symmetric order-3 tensor  $T = \sum_{i=1}^n a_i^{\otimes 3}$ , we say the tensor  $T$  is  $\rho$ -bounded if  $\max_{i \in [n]} \|a_i\|_2 \leq \rho$ .

*Cumulants.* The cumulants of a random vector  $X$  are a sequence of tensors related to the moment tensors of  $X$ :  $K_1(X), K_2(X), K_3(X), \dots$ . We only state the properties we need; see [32] for an introduction. We have  $K_1(X) = \mathbb{E}[X], K_2(X) = \text{cov}(X), K_3(X) = \mathbb{E}[(X - \mathbb{E}[X])^{\otimes 3}]$ . Cumulants have the property that for two independent random variables  $X, Y$  we have  $K_m(X+Y) = K_m(X) + K_m(Y)$ . The first two cumulants of a standard Gaussian random vector are the mean and the covariance matrix; all subsequent cumulants are zero.

*The simultaneous diagonalization algorithm* [28]. The basic idea of the simultaneous diagonalization algorithm<sup>2</sup> to decompose a symmetric order-3 tensor with linearly independent components  $a_1, \dots, a_d \in \mathbb{R}^d$  is the following: for random unit vectors  $x, y \in \mathbb{R}^d$ , compute the (right) eigenvectors of  $T_x T_y^{-1}$ . With probability 1, the set of eigenvectors is equal to the set of directions of  $a_i$ s (the eigenvectors recover the  $a_i$ s up to sign and norm). We use a version that allows for the number of  $a_i$ s to be less than  $d$  and that includes an error analysis [17, 16].

**3. Overcomplete order-3 tensor decomposition.** We consider the problem of decomposing (recovering  $a_i$ s) a symmetric order-3 tensor  $T \in \mathbb{R}^{d \times d \times d}$  of rank  $n$ :

$$(3.1) \quad T = \sum_{i \in [n]} a_i^{\otimes 3}.$$

When the  $a_i$ s are linearly independent, the simultaneous diagonalization algorithm efficiently recovers them, given  $T$ . But it has no guarantees if the components are linearly dependent. Our main idea for the linearly dependent case is that it is still possible that a large subset  $\{a_1, \dots, a_r\}$  of components is linearly independent, so if we cancel out the other components,  $\{a_{r+1}, \dots, a_n\}$ , the residual tensor can be efficiently decomposed via the simultaneous diagonalization algorithm. To cancel the other components, we search for a vector  $x$  orthogonal to them so that  $T_x$  only involves the linearly independent components. A random or grid search for an approximately orthogonal  $x$  is efficient if the number of components to cancel out is small.

For clarity, we now describe an idealized version of our algorithm as if we had two vectors  $x, y$  that are exactly orthogonal to the other components. (The actual algorithm uses a random search to find  $x, y$ .) We also want  $x, y$  to be *generic* with this orthogonality property, so that they can also play the roles of  $x, y$  in the simultaneous diagonalization algorithm (see section 2). Specifically, the genericity here is that the eigenvalues of  $T_x T_y^{-1}$  are distinct. In that case, the eigendecomposition of  $T_x T_y^{-1}$  recovers the directions of  $\{a_1, \dots, a_r\}$ . Then, a linear system of equations provides the lengths of  $\{a_1, \dots, a_r\}$ . Once  $\{a_1, \dots, a_r\}$  is recovered, the components of  $T$  associated to them can be removed from  $T$  (deflation) and the simultaneous diagonalization algorithm can be applied a second time to the residual tensor to recover  $\{a_{r+1}, \dots, a_n\}$ .

**3.1. Approximation algorithm and main theorem.** In the previous discussion, we argued that given  $x, y$ , and  $T$  with exact properties one can decompose  $T$ . In this subsection, we show that by repeatedly trying random choices we can find  $x, y$  nearly orthogonal to  $a_{r+1}, \dots, a_{r+k}$ . In practice, instead of the true tensor  $T$ , we usually have only an approximation  $\tilde{T}$  of it,

---

<sup>2</sup>The simultaneous diagonalization algorithm has been erroneously called Jennrich’s algorithm; see [24] for a discussion.



and to be effective in this situation our algorithm comes with a robustness analysis that shows that if  $\tilde{T}$  is close to  $T$ , then the output is close to the true components of  $T$ . Our formal statements are [Algorithm 2](#) and [Theorem 3.1](#). Our algorithm uses the simultaneous diagonalization algorithm ([Algorithm 1](#), as presented in [\[17\]](#)) as a subroutine.

---

**Algorithm 1.** DIAGONALIZE [\[17\]](#).

---

**Inputs:**  $M_\mu, M_\lambda \in \mathbb{R}^{d \times d}$ , number of vectors  $r$ .

- 1: compute the SVD of  $M_\mu = VDU^\top$ . Let  $W$  be matrix whose columns are the left singular vectors (columns of  $V$ ) corresponding to the top  $r$  singular values;
- 2: compute  $M = (W^\top M_\mu W)(W^\top M_\lambda W)^{-1}$ ;
- 3: compute the eigendecomposition:  $M = P\Lambda P^{-1}$ ;

**Outputs:** columns of  $WP$ .

---



---

**Algorithm 2.** Approximate tensor decomposition.

---

**Inputs:** tensor  $\tilde{T} \in \mathbb{R}^{d \times d \times d}$ , error tolerance  $\varepsilon$ , tensor rank  $n$ , overcompleteness  $k$ , upper bound  $M$  on  $\|a_i\|_2$  for  $i \in [n]$ . Let  $r = n - k$  (Kruskal rank).

- 1: **repeat**
- 2:   pick  $x, y$  i.i.d. uniformly at random in  $\mathcal{S}^{d-1}$ ;
- 3:   invoke [Algorithm 1](#) with  $\tilde{T}_x, \tilde{T}_y$ , and  $r$ . Denote the outputs by  $\tilde{a}_i$  for  $i \in [r]$ ;
- 4:   solve the least squares problem:  $\min_{\xi_1, \dots, \xi_r} \|\tilde{A}_r \text{diag}(\xi_i \langle x, \tilde{a}_i \rangle) \tilde{A}_r^\top - \tilde{T}_x\|_2$ ;
- 5:   set  $R = \tilde{T} - \sum_{i \in [r]} \xi_i \tilde{a}_i^{\otimes 3}$ ;
- 6:   pick  $x', y'$  i.i.d. uniformly at random in  $\mathcal{S}^{d-1}$ ;
- 7:   invoke [Algorithm 1](#) with  $R_{x'}, R_{y'}$ , and  $k$ . Denote the outputs by  $\tilde{a}_{r+i}$  for  $i \in [k]$ ;
- 8:   solve the least squares problem:  $\min_{\xi_{r+1}, \dots, \xi_{r+k}} \|\tilde{A}_{>r} \text{diag}(\xi_{r+i} \langle x', \tilde{a}_{r+i} \rangle) \tilde{A}_{>r}^\top - R_{x'}\|_2$ ;
- 9:   reconstruct the tensor  $T' = \sum_{i \in [r+k]} \xi_i \tilde{a}_i^{\otimes 3}$ ;
- 10: **until**  $\|T' - \tilde{T}\|_F \leq \varepsilon$ ,  $\max_{i \in [r+k]} |\xi_i|^{1/3} \leq 2M$

**Outputs:**  $\tilde{a}_i := \xi_i^{1/3} \tilde{a}_i$  for  $i \in [r+k]$ .

---

**Theorem 3.1 (correctness of [Algorithm 2](#)).** *Let  $T = \sum_{i \in [r+k]} a_i^{\otimes 3}$ ,  $1 \leq k \leq (r-2)/2$ , and  $a_i \in \mathbb{R}^d$ . Let  $A = [a_1, \dots, a_{r+k}]$  and  $\text{K-rank}_\tau(A) \geq r$ . Let  $\tau > 0$ ,  $M \geq \max_{i \in [r+k]} \|a_i\|_2$ ,  $0 < m \leq \min_{i \in [r+k]} \|a_i\|_2$ , and  $0 < \varepsilon_{out} \leq \min\{1, m^3\}$ . There exist polynomials  $\text{poly}_{3.1}(d, \tau, M)$ ,  $\text{poly}'_{3.1}(d, \tau, M, m^{-1})$ , such that if  $\varepsilon_{in} \leq \varepsilon_{out} / \text{poly}'_{3.1}$  and  $\tilde{T}$  is a tensor such that  $\|T - \tilde{T}\|_F \leq \varepsilon_{in}$ , then [Algorithm 2](#) on inputs  $\tilde{T}$  and  $\varepsilon = \varepsilon_{out} / \text{poly}_{3.1}$  outputs vectors  $\tilde{a}_1, \dots, \tilde{a}_{r+k}$  such that for some permutation  $\pi$  of  $[r+k]$ , we have  $\|a_{\pi(i)} - \tilde{a}_i\|_2 \leq \varepsilon_{out}$  for all  $i \in [r+k]$ . The expected running time is at most  $\text{poly}(d^k, \varepsilon_{out}^{-k}, \tau^k, M^k, m^{-k})$ .*

*Proof idea (of [Theorem 3.1](#)).* The proof has three parts. First, we show that if  $T'$  (with which the algorithm finishes) is close to  $\tilde{T}$  and has bounded components, then the components of  $T'$ ,  $\{\tilde{a}_i = \xi_i^{1/3} \tilde{a}_i : i \in [r+k]\}$ , are close to those of  $T$ . In the second part, we show that, assuming good  $x, y, x', y'$  have been found, the algorithm indeed finishes with a tensor  $T'$  that is close to  $\tilde{T}$  (and therefore close to  $T$  via the triangle inequality). We also show how the

error propagates. In the third part, we show the probabilistic bounds that guarantee efficient search of good  $x, y, x', y'$ .

The first part follows from [6, Theorem 2.6] (the version we need is [Theorem 4.1](#) here).

We now informally state what *good*  $x, y$  means. Note that in the idealized case,  $x, y$  are chosen to be orthogonal to  $k = n - r$  vectors and to be generic, meaning that  $T_x T_y^{-1}$  has distinct eigenvalues. In [Algorithm 2](#), we rely on random search to find *good*  $x, y$ , namely  $x, y$  that satisfy the following:

1. nearly orthogonal to last  $k$  terms:  $|\langle y, \hat{a}_{r+i} \rangle|$  are small for  $i \in [k]$  (and similarly for  $x$ );
2. nonorthogonality on first  $r$  terms:  $|\langle y, \hat{a}_i \rangle|$  are lower bounded for  $i \in [r]$  (and similarly for  $x$ );
3. the eigenvalues of  $T_x T_y^{-1}$ ,  $\langle x, \hat{a}_i \rangle / \langle y, \hat{a}_i \rangle$ , are well-separated.

Properties 1 and 2 guarantee that we have  $r$  components with noise after contraction, and property 3 guarantees that the simultaneous diagonalization algorithm can be applied to contracted matrices. We will revisit these properties in [subsection 4.3](#). There are also similar properties for  $x', y'$ .

For the second part, we will assume that we have found good vectors  $x, y$ . [Theorem 4.3](#) (from [16]) and [Lemma 4.4](#) guarantee that we can simultaneously diagonalize matrices  $\tilde{T}_x$  and  $\tilde{T}_y$  using the simultaneous diagonalization algorithm ([Algorithm 1](#)), and the outputs are close to the directions of  $a_i$ s. [Lemma 4.5](#) shows that we can recover approximately the lengths of  $a_i$ s by solving a least squares problem once we have the directions. At this point, we completed the recovery of  $r$  components. [Lemma 4.6](#) shows that when the deflation error is small, the residual tensor  $R$  can be decomposed in the same way and the last  $k$  directions are recovered. At the end of the second part, [Lemma 4.7](#) shows that the lengths of the last  $k$  components are approximately recovered.

The third part is shown in [Lemmas 4.10](#) and [4.11](#). ■

*Remark 3.2.* The constraint  $1 \leq k \leq (r-2)/2$  on the rank is because of Kruskal’s theorem: we need  $2(r+k) + 2 \leq 3r$  to guarantee identifiability.

*Remark 3.3.* [Theorem 3.1](#) has an immediate extension to order- $3p$  symmetric tensors for integer  $p > 1$  by “batching” each set of  $p$  modes together and reshaping into a  $d^p \times d^p \times d^p$  tensor. However, for higher order tensors, additional tools are available. Hence we restrict ourselves to the (in this sense) harder case of order-3 tensors.

**4. Proof of [Theorem 3.1](#).** In this section, we implement the three parts mentioned in the “proof idea” in [subsections 4.1](#) to [4.3](#), respectively. We combine them in [subsection 4.4](#).

**4.1. Uniqueness of decomposition.** We show that if [Algorithm 2](#) satisfies its termination condition, then its outputs are close to the components of  $T$ . We deduce this directly from the following known result on the stability of tensor decompositions.

[Theorem 4.1](#) (see [6, Theorem 5]). *Suppose a rank  $R$  tensor  $T = \sum_{i \in [R]} a_i^{\otimes 3} \in \mathbb{R}^{d \times d \times d}$  is  $\rho$ -bounded. Let  $A = [a_1, \dots, a_R]$  with  $\text{3K-rank}_\tau(A) \geq 2R + 2$ . Then, for every  $\varepsilon' \in (0, 1)$ , there exists  $\varepsilon = \varepsilon' / \text{poly}_{4.1}(R, \tau, \rho, d)$  for a fixed polynomial  $\text{poly}_{4.1}$  so that for any other  $\rho'$ -bounded decomposition  $T' = \sum_{i \in [R]} (a'_i)^{\otimes 3}$  with  $\|T' - T\|_F \leq \varepsilon$ , there exist a permutation matrix  $\Pi$  and diagonal matrix  $\Lambda$  such that  $\|\Lambda^3 - I\|_F \leq \varepsilon'$  and  $\|A' - A\Pi\Lambda\|_F \leq \varepsilon'$ .*



The original statement in [6] explicitly assumes that  $T$  (the sum of  $R$  rank-1 tensors) has rank  $R$ , but this assumption is redundant: a tensor  $T = \sum_{i \in [R]} a_i^{\otimes 3} \in \mathbb{R}^{d \times d \times d}$  with  $3\text{K-rank}(A) \geq 2R + 2$  cannot have another decomposition with less than  $R$  terms because of Kruskal's uniqueness theorem [25, Theorem 4a]. Also, the original statement in [6] is for the nonsymmetric case and we only state here the version we need, specialized to the symmetric case. This restatement is not completely obvious because a symmetric tensor with minimal length symmetric decomposition of length  $R$  (i.e., with symmetric rank equal to  $R$ ) could have a nonsymmetric decomposition of shorter length in general. But under Kruskal's condition,  $3\text{K-rank}_\tau(A) \geq 2R + 2$  (implied by the robust Kruskal condition in Theorem 4.1), Kruskal's uniqueness theorem [25, Theorem 4a] implies that the symmetric and the nonsymmetric decompositions (and ranks) of such a  $T$  coincide.

Note that in Theorem 4.1 a scaling matrix  $\Lambda$  is introduced. We will use the following corollary instead to have a handier result without the scaling matrix.

**Corollary 4.2.** *In the setting of Theorem 4.1, there exists a polynomial  $\text{poly}_{4.2}(R, \tau, \rho, \rho', d)$  such that if  $\varepsilon' \in (0, 1)$  and  $\varepsilon = \varepsilon' / \text{poly}_{4.2}(R, \tau, \rho, \rho', d)$ , then for any other  $\rho'$ -bounded decomposition  $T' = \sum_{i \in [R]} (a'_i)^{\otimes 3}$  with  $\|T' - T\|_F \leq \varepsilon$ , there exists a permutation  $\pi$  of  $[R]$  such that for all  $i \in [R]$ ,  $\|a_{\pi(i)} - a'_i\|_2 \leq \varepsilon'$ .*

*Proof.* We assume that the permutation is the identity. Let  $c = (1 + 4\rho/3)$  and  $\text{poly}_{4.2} = c \text{poly}_{4.1}$ . By Theorem 4.1, we have that for each  $i \in [R]$ :  $\|a'_i - \lambda_i a_i\|_2 \leq c^{-1} \varepsilon'$  and  $|\lambda_i^3 - 1| \leq c^{-1} \varepsilon'$ . Since  $|x - 1| \leq 4|x^3 - 1|/3$  for all  $x \in \mathbb{R}$ , the second inequality implies that  $|\lambda_i - 1| \leq 4|\lambda_i^3 - 1|/3 \leq 4c^{-1} \varepsilon'/3$ . Therefore  $\|a'_i - a_i\|_2 \leq \|a'_i - \lambda_i a_i\|_2 + |\lambda_i - 1| \|a_i\|_2 \leq (1 + 4\rho/3)c^{-1} \varepsilon' = \varepsilon'$ . ■

**4.2. Robust decomposition.** In this subsection, we will derive the forward error propagation of Algorithm 2, i.e., how the output error depends on the input error in each step of Algorithm 2. We will assume throughout this subsection that we already have two unit vectors  $x, y$  that are nearly orthogonal to  $\hat{a}_{r+1}, \dots, \hat{a}_{r+k}$ , that is,  $|\langle x, \hat{a}_{r+i} \rangle|, |\langle y, \hat{a}_{r+i} \rangle| \leq \theta$  for  $i \in [k]$ , where  $\theta$  will be chosen later, and  $\text{K-rank}_\tau(A) \geq r$ . Let  $E_{\text{in}} = T - \hat{T}$  be the input error tensor. Also recall that  $\|a_i\| \in [m, M]$ . We summarize the roadmap of this subsection in Figure 1.

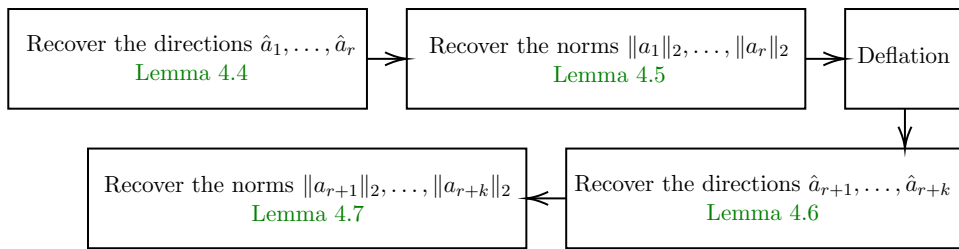


Figure 1. Roadmap of subsection 4.2.

**Part 1: Robust diagonalization.** We first cite the robust analysis of Algorithm 1.

**Theorem 4.3** (see [16, Theorem 5.4, Lemmas 5.1, 5.2]). *Let  $T_\mu = \sum_{i \in [r]} \mu_i a_i a_i^\top = A \text{diag}(\mu) A^\top$ ,  $T_\lambda = \sum_{i \in [r]} \lambda_i a_i a_i^\top = A \text{diag}(\lambda) A^\top$ ,  $A = [a_1, \dots, a_r]$ ,  $a_i \in \mathbb{R}^d$ ,  $\|a_i\| = 1$ ,*

$\lambda_i, \mu_i \in \mathbb{R}$  for  $i \in [r]$ . Suppose **(1)**  $\sigma_r(A) > 0$ , **(2)**  $(\forall i) 0 < k_l \leq |\mu_i|, |\lambda_i| \leq k_u$ , and **(3)**  $(\forall i \neq j) |\mu_i/\lambda_i - \mu_j/\lambda_j| \geq \alpha > 0$ . Let  $0 < \varepsilon_{4.3} < 1$ , and let  $\tilde{T}_\mu, \tilde{T}_\lambda$  be matrices such that  $\|T_\mu - \tilde{T}_\mu\|_F, \|T_\lambda - \tilde{T}_\lambda\|_F \leq \frac{\varepsilon_{4.3} k_l^2 \sigma_r(A)^3 \min\{\alpha, 1\}}{2^{11} \kappa(A) k_u r^2}$ . Then [Algorithm 1](#) on input  $\tilde{T}_\mu, \tilde{T}_\lambda$  outputs unit vectors  $\tilde{a}_1, \dots, \tilde{a}_r$  such that for some permutation  $\pi$  of  $[r]$  and signs  $s_1, \dots, s_r \in \{\pm 1\}$  and for all  $i \in [r]$  we have  $\|a_{\pi(i)} - s_i \tilde{a}_i\| \leq \varepsilon_{4.3}$ . It runs in time  $\text{poly}(d, 1/\alpha, 1/k_l, 1/\sigma_r(A_r), 1/\varepsilon_{4.3})$ .

Now we apply [Theorem 4.3](#) to our case: let  $E_x = T_x - \tilde{T}_x$  and  $E_y = T_y - \tilde{T}_y$ . Write  $\tilde{T}_x = \hat{A}_r D_x \hat{A}_r^\top + \hat{A}_{>r} D'_x \hat{A}_{>r}^\top + (E_{\text{in}})_x$ , where  $\hat{A}_r$  contains  $\hat{a}_i$ s as columns,  $D_x = \text{diag}(\|a_i\|^3 \langle x, \hat{a}_i \rangle)$  for  $i \in [r]$ , and  $\hat{A}_{>r}$  contains  $\hat{a}_{r+i}$ s,  $D'_x = \text{diag}(\|a_{r+i}\|^3 \langle x, \hat{a}_{r+i} \rangle)$  for  $i \in [k]$ . Then we have

$$(4.1) \quad \|E_x\|_F = \|\hat{A}_{>r} D'_x \hat{A}_{>r}^\top + (E_{\text{in}})_x\|_F \leq kM^3\theta + \varepsilon_{\text{in}},$$

and similarly for  $E_y$ . The following lemma guarantees the correctness of step 3 in [Algorithm 2](#).

**Lemma 4.4 (direction estimation).** Let  $\tilde{a}_1, \dots, \tilde{a}_r$  be the outputs of step 3 in [Algorithm 2](#). If **(1)**  $\forall i \in [r]: 0 < k_l/m^3 \leq |\langle x, \hat{a}_i \rangle|, |\langle y, \hat{a}_i \rangle| \leq 1$ , and **(2)**  $\forall i, j \in [r], i \neq j: |\langle x, \hat{a}_i \rangle / \langle y, \hat{a}_i \rangle - \langle x, \hat{a}_j \rangle / \langle y, \hat{a}_j \rangle| \geq \alpha > 0$ , then there exist signs  $s_1, \dots, s_r \in \{\pm 1\}$  and a permutation  $\pi$  of  $[r]$  such that for all  $i \in [r]$ ,

$$\|\hat{a}_{\pi(i)} - s_i \tilde{a}_i\| \leq \varepsilon_{4.4} := \frac{2^{11} \tau^4 M^7 r^{5/2} (kM^3\theta + \varepsilon_{\text{in}})}{k_l^2 \min\{\alpha, 1\}}.$$

This step runs in time  $\text{poly}(d, \alpha^{-1}, k_l^{-1}, \tau, M, \varepsilon_{4.4}^{-1})$ .

*Proof.* Condition 1 in [Theorem 4.3](#) holds since  $\text{K-rank}_\tau(A) \geq r: \sigma_r(\hat{A}_r) \geq \sigma_r(A_r)/M \geq 1/(\tau M)$ . Conditions 2 and 3 in [Theorem 4.3](#) hold because of our assumptions. Combining (4.1) and  $\text{K-rank}_\tau(A) \geq r$ , which implies  $\sigma_r(\hat{A}_r)^3 \kappa(\hat{A}_r)^{-1} = \sigma_r(\hat{A}_r)^4 \sigma_1(\hat{A}_r)^{-1} \geq (\sqrt{r} \tau^4 M^4)^{-1}$ , the assumptions of [Theorem 4.3](#) are satisfied with parameter  $k_u = M^3$ . The claim follows. ■

Since  $x, y$  are actually chosen at random, we provide the probability for assumptions of [Lemma 4.4](#) to hold in [subsection 4.3](#).

**Part 2: Norm estimation.** The next step is to recover  $\|a_i\|_2$ . This can be done by solving the least squares problem in step 4. To see this, one can verify that when  $\tilde{a}_i = \hat{a}_i$  and  $\tilde{T}_x = T_x$  (no error in earlier steps),  $\xi_i = \|a_i\|_2^3$  is a zero error solution to step 4. The following lemma guarantees that we can approximate the norm via step 4.

**Lemma 4.5 (norm estimation).** Let  $\tilde{b}_1, \dots, \tilde{b}_r$  be the columns of  $(\tilde{A}_r^\dagger)^\top$ . If [Lemma 4.4](#) holds with  $\varepsilon_{4.4} \leq \min\{k_l/(2m^3), (2\sqrt{r}\tau M)^{-1}\}$ , then  $\xi_i = \tilde{T}(x, \tilde{b}_i, \tilde{b}_i) / \langle x, \tilde{a}_i \rangle$  for  $i \in [r]$  is the unique solution to step 4 in [Algorithm 2](#) and for the permutation  $\pi$ , signs  $s_i$  in [Lemma 4.4](#), and all  $i \in [r]$  we have

$$\left| \|a_{\pi(i)}\|_2^3 - s_i \xi_i \right| \leq \varepsilon_{4.5} := 2k_l^{-1} m^3 M^2 [3M\varepsilon_{4.4} + rM\varepsilon_{4.4}^2 + 4\tau^2 (kM^3\theta + \varepsilon_{\text{in}})].$$

*Proof.* For simplicity, we assume the permutation is the identity. We start by showing that  $\sigma_r(\tilde{A}_r) > 0$ , which implies  $\tilde{A}_r^\dagger \tilde{A}_r = I_r$  and thus  $\tilde{b}_i$  is orthogonal to  $\tilde{a}_j$  for  $i, j \in [r], i \neq j$ . By [Lemma 4.4](#), the distance between corresponding columns of  $\tilde{A}_r \text{diag}(s_i)$  and  $\hat{A}_r$  is at most  $\varepsilon_{4.4}$ . Therefore by [Theorem B.1](#) we have  $|\sigma_r(\tilde{A}_r \text{diag}(s_i)) - \sigma_r(\hat{A}_r)| \leq \|\tilde{A}_r \text{diag}(s_i) - \hat{A}_r\|_2 \leq \sqrt{r}\varepsilon_{4.4}$ , which implies

$$(4.2) \quad \sigma_r(\tilde{A}_r) = \sigma_r(\tilde{A}_r \text{diag}(s_i)) \geq \sigma_r(\hat{A}_r) - \sqrt{r}\varepsilon_{4.4} \geq (\tau M)^{-1} - \sqrt{r}\varepsilon_{4.4} \geq 1/(2\tau M).$$

Next, we show that  $\xi_i$  is the unique solution to step 4. We restate the least squares problem in a matrix-vector product form,  $\min_{\xi_i} \|\tilde{A}^{\odot 2}[\langle x, \tilde{a}_1 \rangle \xi_1, \dots, \langle x, \tilde{a}_r \rangle \xi_r]^\top - \text{vec}(\tilde{T}_x)\|_2$ , where  $\tilde{A}^{\odot 2} = [\text{vec}(\tilde{a}_1 \tilde{a}_1^\top), \dots, \text{vec}(\tilde{a}_r \tilde{a}_r^\top)] \in \mathbb{R}^{d^2 \times r}$ . It follows that  $\sigma_r(\tilde{A}^{\odot 2}) = \sigma_r(\tilde{A}_r)^2 > 0$  and thus the solution is unique. Let  $\tilde{B}^{\odot 2} = [\text{vec}(\tilde{b}_1 \tilde{b}_1^\top), \dots, \text{vec}(\tilde{b}_r \tilde{b}_r^\top)]^\top$ , and notice that  $\tilde{B}^{\odot 2} \tilde{A}^{\odot 2} = I_r$ . The solution to the least squares problem is then given by  $[\langle x, \tilde{a}_1 \rangle \xi_1, \dots, \langle x, \tilde{a}_r \rangle \xi_r]^\top = \tilde{B}^{\odot 2} \text{vec}(\tilde{T}_x) = [\tilde{b}_1^\top \tilde{T}_x \tilde{b}_1, \dots, \tilde{b}_r^\top \tilde{T}_x \tilde{b}_r]^\top$ , which implies  $\xi_i = \tilde{T}(x, \tilde{b}_i, \tilde{b}_i) / \langle x, \tilde{a}_i \rangle$ .

Finally, we show that  $s_i \xi_i$  is close to  $\|a_i\|_2^3$ . The deviation of  $s_i \xi_i$  from  $\|a_i\|_2^3$  is bounded by

$$(4.3) \quad \begin{aligned} & \left| \|a_i\|_2^3 - s_i \xi_i \right| = \left| \|a_i\|_2^3 - \frac{1}{\langle x, s_i \tilde{a}_i \rangle} \left( \sum_{j \in [r]} \langle x, a_j \rangle \langle \tilde{b}_i, a_j \rangle^2 + \tilde{b}_j^\top E_x \tilde{b}_j \right) \right| \\ & \leq \underbrace{\left| \frac{\langle x, \hat{a}_i \rangle \langle \tilde{b}_i, \hat{a}_i \rangle^2}{\langle x, s_i \tilde{a}_i \rangle} - 1 \right|}_{\substack{\text{small when } s_i \tilde{a}_i \text{ close to } \hat{a}_i \\ \text{(note that } \langle \tilde{b}_i, \tilde{a}_i \rangle = 1)}} \|a_i\|_2^3 + \sum_{j \in [r], j \neq i} \left( \|a_j\|_2^3 \underbrace{\left| \frac{\langle x, \hat{a}_j \rangle \langle \tilde{b}_i, \hat{a}_j \rangle^2}{\langle x, s_i \tilde{a}_i \rangle} \right|}_{\substack{\text{small when } s_j \tilde{a}_j \text{ close to } \hat{a}_j \\ \text{(note that } \langle \tilde{b}_i, \tilde{a}_j \rangle = 0)}} + \underbrace{\left| \frac{\tilde{b}_i^\top E_x \tilde{b}_i}{\langle x, s_i \tilde{a}_i \rangle} \right|}_{\text{error from } E_x} \right). \end{aligned}$$

We analyze the deviation of each term in (4.3). By standard arguments, using the triangle and Cauchy–Schwarz inequalities, we have for all  $i, j \in [r]$ ,

$$(4.4) \quad \begin{aligned} & |\langle x, s_i \tilde{a}_i \rangle| \geq |\langle x, \hat{a}_i \rangle| - \varepsilon_{4.4} \geq k_l/m^3 - \varepsilon_{4.4} \geq k_l/(2m^3), \\ & |\langle x, s_j \tilde{a}_j \rangle - \langle x, \hat{a}_j \rangle| \leq \varepsilon_{4.4}, \quad |\langle \tilde{b}_i, s_j \tilde{a}_j \rangle - \langle \tilde{b}_i, \hat{a}_j \rangle| \leq \varepsilon_{4.4}, \end{aligned}$$

where the first line comes from the assumptions of the lemma, and the last line follows from Lemma 4.4. Notice that  $\tilde{b}_i$  is orthogonal to  $\tilde{a}_j$  for  $j \neq i$ , and  $\langle \tilde{b}_i, \tilde{a}_i \rangle = 1$ . Equation (4.4) implies that

$$(4.5) \quad \left| \frac{\langle x, \hat{a}_i \rangle \langle \tilde{b}_i, \hat{a}_i \rangle^2}{\langle x, s_i \tilde{a}_i \rangle} - 1 \right| \leq 6k_l^{-1} m^3 \varepsilon_{4.4}, \quad \left| \frac{\langle x, \hat{a}_j \rangle \langle \tilde{b}_i, \hat{a}_j \rangle^2}{\langle x, s_i \tilde{a}_i \rangle} \right| \leq 2k_l^{-1} m^3 \varepsilon_{4.4}^2.$$

The last term in (4.3) is bounded by

$$(4.6) \quad \left| \frac{\tilde{b}_i^\top E_x \tilde{b}_i}{\langle x, s_i \tilde{a}_i \rangle} \right| \leq 2k_l^{-1} m^3 \|E_x\|_2 \|\tilde{b}_i\|_2^2 \leq 2k_l^{-1} m^3 \|E_x\|_F \sigma_r(\tilde{A}_r)^{-2} \leq 8k_l^{-1} m^3 \tau^2 M^2 \|E_x\|_F,$$

where the second inequality follows from the definition of  $\tilde{b}_i$ , and the last inequality applies (4.2). Combining (4.1), (4.3), (4.5), and (4.6) gives the desired result.  $\blacksquare$

**Part 3: Deflation.** After we deflate  $T$  with the previously recovered  $r$  components, the induced error with respect to the exact deflation  $\sum_{i=r+1}^{r+k} a_i^{\otimes 3}$  is given by  $E' = E_{\text{in}} + \sum_{i \in [r]} (a_i^{\otimes 3} - \xi_i \tilde{a}_i^{\otimes 3})$ . Now we show that the remaining tensor can be decomposed with the same strategy via step 7 in Algorithm 2.

**Lemma 4.6 (direction estimation).** *Let  $\tilde{a}_{r+1}, \dots, \tilde{a}_{r+k}$  be the outputs of step 7 in Algorithm 2. If (1)  $\forall i \in [k]: 0 < k'_l/m^3 \leq |\langle x', \hat{a}_{r+i} \rangle|, |\langle y', \hat{a}_{r+i} \rangle| \leq 1$ , and (2)  $\forall i, j \in [k], i \neq j$ :*

$|\langle x', \hat{a}_{r+i} \rangle / \langle y', \hat{a}_{r+i} \rangle - \langle x', \hat{a}_{r+j} \rangle / \langle y', \hat{a}_{r+j} \rangle| \geq \alpha' > 0$ , then there exist signs  $s_{r+1}, \dots, s_{r+k} \in \{\pm 1\}$  and a permutation  $\pi'$  of  $[k]$  such that for all  $i \in [k]$ ,

$$\|\hat{a}_{r+\pi'(i)} - s_{r+i} \tilde{a}_{r+i}\|_2 \leq \varepsilon_{4.6} := \frac{2^{11} \tau^4 M^7 k^{5/2} \|E'\|_F}{(k'_l)^2 \min\{\alpha', 1\}}.$$

This step runs in time  $\text{poly}(d, k'_l{}^{-1}, \alpha'^{-1}, \tau, M, \varepsilon_{4.6}^{-1})$ .

*Proof.* The proof is similar to the proof of Lemma 4.4 and is thus omitted here. ■

With  $\tilde{a}_{r+1}, \dots, \tilde{a}_{r+k}$ , we can further approximate the norm of  $a_{r+1}, \dots, a_{r+k}$ , in the same way we did for the first  $r$  components, via step 8. The following lemma guarantees it works.

**Lemma 4.7 (norm estimation).** *Let  $\tilde{b}_{r+1}, \dots, \tilde{b}_{r+k}$  be the columns of  $(\tilde{A}_{\geq r}^\dagger)^\top$ . If Lemma 4.6 holds with  $\varepsilon_{4.6} \leq \min\{k'_l / (2m^3), (2\sqrt{k}\tau M)^{-1}\}$ , then  $\xi_{r+i} = R(x', \tilde{b}_{r+i}, \tilde{b}_{r+i}) / \langle x', \tilde{a}_{r+i} \rangle$ , for  $i \in [k]$  is the unique solution to step 8 in Algorithm 2 and for the permutation  $\pi'$ , signs  $s_{r+i}$  in Lemma 4.6, and all  $i \in [k]$  we have*

$$\| \|a_{r+\pi'(i)}\|^3 - s_{r+i} \xi_{r+i} \| \leq \varepsilon_{4.7} := 2k'_l{}^{-1} m^3 M^2 [3M\varepsilon_{4.6} + kM\varepsilon_{4.6}^2 + 4\tau^2 \|E'\|_F].$$

*Proof.* The proof is similar to the proof of Lemma 4.5 and is thus omitted here. ■

**4.3. Probability bounds.** We give here bounds on the probability of finding good  $x, y, x', y'$  so that Algorithm 2 succeeds with positive probability. Throughout this subsection, let  $x, y$  be two i.i.d. random vectors distributed uniformly on  $\mathcal{S}^{d-1}$ , and let  $a_1, a_2, \dots, a_{r+k}$  be such that  $\|a_i\| \in [m, M]$  and  $\text{K-rank}_\tau([a_1, \dots, a_{r+k}]) \geq r$ , which implies that their directions satisfy  $\text{K-rank}_{\tau M}([\hat{a}_1, \dots, \hat{a}_{r+k}]) \geq r$ .

We first list the events for good  $x, y$  to hold to apply Lemma 4.4:

1. nearly orthogonal to last  $k$  terms:  $\mathcal{E}_{1,y} = \{\forall i \in [k], |\langle y, \hat{a}_{r+i} \rangle| \leq \theta\}$ ;
2. nonorthogonality on first  $r$  terms:  $\mathcal{E}_{2,y} = \{\forall i \in [r], |\langle y, \hat{a}_i \rangle| \geq k_l / m^3\}$ ;
3. the eigenvalue gap:  $\mathcal{E}_3 = \{\forall i \neq j, i, j \in [r], |\langle x, \hat{a}_i \rangle / \langle y, \hat{a}_i \rangle - \langle x, \hat{a}_j \rangle / \langle y, \hat{a}_j \rangle| \geq \alpha > 0\}$ .

We have similar events  $\mathcal{E}_{1,x}, \mathcal{E}_{2,x}$ . Note that in this subsection  $k_l, \theta$ , and  $\alpha$  are considered as fixed parameters.

The structure of this subsection is stated as follows: we will first demonstrate our proof idea for controlling the probability of  $\mathcal{E}_{1,y} \cap \mathcal{E}_{2,y}$ . After presenting our idea, we will first analyze the probability of  $\mathcal{E}_{1,y} \cap \mathcal{E}_{2,y}$  and then the probability of  $\mathcal{E}_{1,x} \cap \mathcal{E}_{2,x} \cap \mathcal{E}_3$  when conditioned on the other events of  $y$ . Finally, we will collect these subevents and give the probability that all of them will hold.

It seems that direct union bound-type arguments are insufficient and some nontrivial conditioning is necessary: First,  $\mathcal{E}_{1,x}$  happens with small probability, as meaningful values of  $\theta$  have to be much smaller than  $k_l / m^3$  and  $\alpha$ . Naively applying the union bound on some events and analyzing the rest does not give enough wiggle room for a positive probability. Besides,  $\mathcal{E}_3$  is the most complicated in the sense that it controls the difference of two ratios. As we will see later, after conditioning on  $y$ , the ideas of analyzing  $\mathcal{E}_{1,y}, \mathcal{E}_{2,y}$  can be reused for the rest of the events, which makes the analysis easier to follow.

We now state the idea of our argument to bound the probability of  $\mathcal{E}_{1,y} \cap \mathcal{E}_{2,y}$ .

*Bands argument.* We analyze the events geometrically and replace random unit vectors by random Gaussian vectors together with concentration of their norm. Let  $z$  be a random Gaussian vector, and let  $a$  and  $b$  be two unit vectors. An event of the form  $\{|\langle z, a \rangle| \leq t_1\}$  corresponds to a band, while an event like  $\{|\langle z, b \rangle| \geq t_2\}$  corresponds to the complement of a band. We call them bands of type I and type II, denoted by  $\mathcal{B}_1$  and  $\mathcal{B}_2$ , respectively. To better illustrate this, we give a demonstration of bands as the shaded areas in Figure 2.



Figure 2. Examples of bands.

The intersection of bands of type I can be lower-bounded with a direct use of the Gaussian correlation inequality, Lemma B.5, while the intersection of bands of different types needs special care. Consider  $\mathcal{B}_1 \cap \mathcal{B}_2$ : when  $\langle a, b \rangle = 0$ , the intersection becomes  $\mathcal{B}_1$  with a rectangular region excluded. In this case, the two bands will be orthogonal, and the two events are independent. In the general case, the excluded region is a parallelogram depending on  $\langle a, b \rangle$ . See Figure 3 for illustration. In the extreme case, two bands are parallel, and hence the probability will be zero when  $t_1 \leq t_2$ . But when  $\langle a, b \rangle$  is not too close to one, we can, when bounding the probability, replace the parallelogram by a slightly larger rectangular region without decreasing the final probability too much, which is shown by the white dashed lines in Figure 3b. This is essentially done by projecting  $b$  onto  $\text{span}\{a\}$  and  $\text{span}\{a\}^\perp$ .



Figure 3. Intersection of bands.

We see that events  $\mathcal{E}_{1,y}, \mathcal{E}_{2,y}$  are the intersection of bands and their probability is the probability measure of their intersection. Specifically, we have  $\mathcal{E}_{1,y} = \cap_{i=1}^k \mathcal{B}_{1,i}$ ,  $\mathcal{E}_{2,y} = \cap_{j=1}^r \mathcal{B}_{2,j}$ , where  $\mathcal{B}_{1,i} := \{|\langle y, \hat{a}_{r+i} \rangle| \leq \theta\}$  and  $\mathcal{B}_{2,j} := \{|\langle y, \hat{a}_j \rangle| \geq k_l/m^3\}$ . For the rest of this subsection, let  $S = \text{span}\{\hat{a}_{r+1}, \dots, \hat{a}_{r+k}\}^\perp$ , let  $S^\perp = \text{span}\{\hat{a}_{r+1}, \dots, \hat{a}_{r+k}\}$ , let  $\text{proj}_S$  be the orthogonal projection onto  $S$ , and let  $\text{proj}_{S^\perp} = I - \text{proj}_S$ . Now we can bound the probability of  $\mathcal{E}_{1,y} \cap \mathcal{E}_{2,y}$ .

**Lemma 4.8.** *If  $k_l > 0$  and  $0 < \theta \leq 2/\sqrt{d}$ , then  $\mathbb{P}[\mathcal{E}_{1,y} \cap \mathcal{E}_{2,y}] \geq p_1 := (\theta\sqrt{d}/8)^k (1/4 - r\sqrt{d}/2\pi\tau M(4k_l/m^3 + \sqrt{k}\tau M\theta))$ .*

*Proof.* Write  $y = z/\|z\|_2$ , where  $z$  is a standard Gaussian random vector. Consider the following events corresponding to  $z$  for  $R_1, R_2$  to be chosen later:  $\mathcal{B}'_{1,i} := \{|\langle z, \hat{a}_{r+i} \rangle| \leq R_1\theta\}$

and  $\mathcal{B}'_{2,j} := \{|\langle z, \hat{a}_j \rangle| \geq R_2 k_l / m^3\}$ . We have

$$\begin{aligned} \mathcal{E}_{1,y} \cap \mathcal{E}_{2,y} &= (\cap_i \mathcal{B}_{1,i}) \cap (\cap_j \mathcal{B}_{2,j}) = (\cap_i \mathcal{B}_{1,i}) \setminus (\cup_j \mathcal{B}_{2,j}^c) \\ &\supseteq \underbrace{(\cap_i \mathcal{B}'_{1,i} \setminus \{\|z\|_2 \leq R_1\})}_{z \text{ nearly orthogonal to } \hat{a}_{r+i} \text{ while } \|z\| > R_1} \setminus \underbrace{(\cup_j ((\mathcal{B}'_{2,j})^c \cup \{\|z\|_2 \geq R_2\}))}_{z \text{ nearly orthogonal to } \hat{a}_j \text{ for some } j \text{ or } \|z\| \geq R_2}. \end{aligned}$$

Set  $\mathcal{E} = \cap_{i \in [k]} \mathcal{B}'_{1,i}$ . Since  $\mathcal{E} \setminus \{\|z\|_2 \leq R_1\} = \mathcal{E} \setminus (\{\|z\|_2 \leq R_1\} \cap \mathcal{E}) \supseteq \mathcal{E} \setminus (\{\|\text{proj}_S z\|_2 \leq R_1\} \cap \mathcal{E})$ ,

$$\begin{aligned} \mathcal{E}_{1,y} \cap \mathcal{E}_{2,y} &\supseteq (\mathcal{E} \setminus (\{\|\text{proj}_S z\|_2 \leq R_1\} \cap \mathcal{E})) \setminus \cup_j ((\mathcal{B}'_{2,j})^c \cup \{\|z\|_2 \geq R_2\}) \\ &= \mathcal{E} \setminus \left( (\{\|\text{proj}_S z\|_2 \leq R_1\} \cap \mathcal{E}) \cup \cup_{j \in [r]} ((\mathcal{B}'_{2,j})^c \cap \mathcal{E}) \cup (\{\|z\|_2 \geq R_2\} \cap \mathcal{E}) \right), \end{aligned}$$

which implies

$$(4.7) \quad \begin{aligned} &\mathbb{P}[\mathcal{E}_{1,y} \cap \mathcal{E}_{2,y}] \\ &\geq \mathbb{P}[\mathcal{E}] - \sum_{j \in [r]} \mathbb{P}[(\mathcal{B}'_{2,j})^c, \mathcal{E}] - \mathbb{P}[\{\|\text{proj}_S z\|_2 \leq R_1\}, \mathcal{E}] - \mathbb{P}[\{\|z\|_2 \geq R_2\}, \mathcal{E}]. \end{aligned}$$

We now bound the terms in (4.7). First,

$$\mathbb{P}[(\mathcal{B}'_{2,j})^c, \mathcal{E}] = \mathbb{P}[|\langle z, \hat{a}_j \rangle| \leq R_2 k_l / m^3 \mid \mathcal{E}] \mathbb{P}[\mathcal{E}].$$

Notice that when conditioning on the event  $|\langle z, \hat{a}_{r+i} \rangle| \leq R_1 \theta$  for  $i \in [k]$  we have

$$(4.8) \quad |\langle z, \text{proj}_{S^\perp} \hat{a}_j \rangle| = |z^\top \hat{A}_{>r} \hat{A}_{>r}^\dagger \hat{a}_j| \leq R_1 \sqrt{k} \theta \|\hat{A}_{>r}^\dagger \hat{a}_j\|_2 \leq R_1 \sqrt{k} \tau M \theta,$$

where the first equality comes from the definition of the projection, the second inequality follows from the conditioning, and the last one comes from the robust Kruskal rank condition. Furthermore, we notice that  $\text{proj}_S \hat{a}_j$  is orthogonal to  $\hat{a}_{r+1}, \dots, \hat{a}_{r+k}$  and the conditioning can therefore be dropped after applying (4.8):

$$\begin{aligned} \mathbb{P}[|\langle z, \hat{a}_j \rangle| \leq R_2 k_l / m^3 \mid \mathcal{E}] &\leq \mathbb{P}[|\langle z, \text{proj}_S \hat{a}_j \rangle| \leq R_2 k_l / m^3 + |\langle z, \text{proj}_{S^\perp} \hat{a}_j \rangle| \mid \mathcal{E}] \\ &\leq \mathbb{P}[|\langle z, \text{proj}_S \hat{a}_j \rangle| \leq R_2 k_l / m^3 + R_1 \sqrt{k} \tau M \theta] \\ &\leq 2(\sqrt{2\pi} \|\text{proj}_S \hat{a}_1\|_2)^{-1} (R_2 k_l / m^3 + R_1 \sqrt{k} \tau M \theta) \\ &\leq \sqrt{2/\pi} \tau M (R_2 k_l / m^3 + R_1 \sqrt{k} \tau M \theta), \end{aligned}$$

where the last two steps follow from bounding the density of a Gaussian distribution from above and the fact that  $\{\hat{a}_j, \hat{a}_{r+1}, \dots, \hat{a}_{r+k}\}$  also satisfies the robust Kruskal rank condition so that  $\|\text{proj}_S \hat{a}_j\|_2 \geq (\tau M)^{-1}$ .

We use the following bounds for the rest of the terms in (4.7):

$$\begin{aligned} \mathbb{P}[\mathcal{E}] &\geq (R_1 \theta / 4)^k \quad (\text{Lemma B.5}), \\ \mathbb{P}[\|\text{proj}_S z\|_2 \leq R_1, \mathcal{E}] &= \mathbb{P}[\|\text{proj}_S z\|_2 \leq R_1] \mathbb{P}[\mathcal{E}] \leq \mathbb{P}[\mathcal{E}] / 2 \quad (\text{set } R_1 = \sqrt{d}/2), \\ \mathbb{P}[\|z\|_2 \geq R_2, \mathcal{E}] &= \mathbb{P}[\|z\|_2 \geq R_2 \mid \mathcal{E}] \mathbb{P}[\mathcal{E}] = (1 - \mathbb{P}[\|z\|_2 \leq R_2 \mid \mathcal{E}]) \mathbb{P}[\mathcal{E}] \\ &\leq (1 - \mathbb{P}[\|z\|_2 \leq R_2]) \mathbb{P}[\mathcal{E}] \quad (\text{Gaussian correlation ineq., Theorem B.4}) \\ &\leq \mathbb{P}[\mathcal{E}] / 4 \quad (\text{Markov's inequality, set } R_2 = 2\sqrt{d}). \end{aligned}$$



Combining the previous estimates, we get  $\mathbb{P}[\mathcal{E}_{1,y} \cap \mathcal{E}_{2,y}] \geq \mathbb{P}[\mathcal{E}](1 - 1/2 - r\sqrt{2/\pi}\tau M(R_2 k_l/m^3 + R_1\sqrt{k}\tau M\theta) - 1/4)$ . The claim follows.  $\blacksquare$

At this point, we are ready to analyze the probability of  $\mathcal{E}_3$ .

**Lemma 4.9.** *In the setting of Lemma 4.8, let  $p_2 = p_1 - (\theta\sqrt{d}/8)^k r^2 \tau M(\sqrt{dk}\theta\tau M k_l^{-1} m^3 + \alpha)$ . Then  $\mathbb{P}[\mathcal{E}_3 \cap \mathcal{E}_{1,x} \cap \mathcal{E}_{2,x} \mid \mathcal{E}_{1,y}, \mathcal{E}_{2,y}] \geq p_2$ .*

*Proof.* We start with our idea to bound the probability of the ‘‘eigenvalue gap’’  $\left| \frac{\langle x, \hat{a}_s \rangle}{\langle y, \hat{a}_s \rangle} - \frac{\langle x, \hat{a}_t \rangle}{\langle y, \hat{a}_t \rangle} \right| \geq \alpha$  for  $s, t \in [r], s \neq t$ . Since we condition on  $|\langle y, \hat{a}_i \rangle|$  not being too small for all  $i \in [r]$ , when further conditioned on  $y$ , we have

$$\mathbb{P}\left[\left| \frac{\langle x, \hat{a}_s \rangle}{\langle y, \hat{a}_s \rangle} - \frac{\langle x, \hat{a}_t \rangle}{\langle y, \hat{a}_t \rangle} \right| \geq \alpha \mid \mathcal{E}_{1,y}, \mathcal{E}_{2,y}\right] = \mathbb{E}\left[\mathbb{P}\left[\underbrace{\left| \frac{\langle x, \hat{a}_s \rangle}{\langle y, \hat{a}_s \rangle} - \frac{\langle x, \hat{a}_t \rangle}{\langle y, \hat{a}_t \rangle} \right|}_{\text{denominators are fixed}} \geq \alpha \mid y\right] \mid \mathcal{E}_{1,y}, \mathcal{E}_{2,y}\right].$$

Therefore it is enough to show a uniform lower bound for  $\mathbb{P}[|\langle x, C_s \hat{a}_s - C_t \hat{a}_t \rangle| \geq \alpha]$ , where  $|C_s|, |C_t|$  are in  $[1, k_l^{-1} m^3]$ . We notice that the set  $\{|\langle x, C_s \hat{a}_s - C_t \hat{a}_t \rangle| \geq \alpha\}$  generates a type II band, denoted by  $\mathcal{B}_{3,st}$ . Therefore the target event is the intersection of  $k$  type I bands  $\mathcal{B}_{1,i}$ ,  $r$  type II bands  $\mathcal{B}_{2,j}$ , and  $\binom{r}{2}$  type II bands  $\mathcal{B}_{3,st}$ . More precisely,

$$\mathbb{P}[\mathcal{E}_3, \mathcal{E}_{1,x}, \mathcal{E}_{2,x} \mid \mathcal{E}_{1,y}, \mathcal{E}_{2,y}] \geq \inf_{|C_s|, |C_t| \in [1, k_l^{-1} m^3]} \mathbb{P}[\cap_{i \in [k]} \mathcal{B}_{1,i}, \cap_{j \in [r]} \mathcal{B}_{2,j}, \cap_{s,t \in [r], s \neq t} \mathcal{B}_{3,st}].$$

We reuse ideas from the proof of Lemma 4.8. Write  $x = u/\|u\|_2$  with  $u$  being standard Gaussian. Consider the following events for  $u$ :  $\mathcal{B}'_{1,i} := \{|\langle u, \hat{a}_{r+i} \rangle| \leq \sqrt{d}\theta/2\}$ ,  $\mathcal{B}'_{2,j} := \{|\langle u, \hat{a}_j \rangle| \geq 2\sqrt{dk_l}/m^3\}$ , and  $\mathcal{B}'_{3,st} := \{|\langle u, C_s \hat{a}_s - C_t \hat{a}_t \rangle| \geq 2\sqrt{d}\alpha\}$ . Set  $\mathcal{E} = \cap_{i \in [k]} \mathcal{B}'_{1,i}$ . With the concentration of  $\|u\|_2$  in  $[\sqrt{d}/2, 2\sqrt{d}]$ , the target probability becomes

$$\begin{aligned} \mathbb{P}[\cap_{i \in [k]} \mathcal{B}_{1,i}, \cap_{j \in [r]} \mathcal{B}_{2,j}, \cap_{s,t \in [r], s \neq t} \mathcal{B}_{3,st}] &\geq \mathbb{P}\left[\mathcal{E} \setminus \left(\left(\{\| \text{proj}_S u \|_2 \leq \sqrt{d}/2\} \cap \mathcal{E}\right) \right. \right. \\ &\quad \left. \left. \cup \cup_{j \in [r]} ((\mathcal{B}'_{2,j})^c \cap \mathcal{E}) \cup (\{\|u\|_2 \geq 2\sqrt{d}\} \cap \mathcal{E}) \cup \cup_{s \neq t \in [r]} ((\mathcal{B}'_{3,st})^c \cap \mathcal{E})\right)\right] \\ (4.9) \quad &\geq p_1 - \sum_{s,t \in [r], s \neq t} \mathbb{P}[\mathcal{E}, (\mathcal{B}'_{3,st})^c]. \end{aligned}$$

Now we consider the summand, which is the intersection of  $k+1$  type I bands. Take  $s = 1, t = 2$  (the rest is similar), and write  $v = C_1 \hat{a}_1 - C_2 \hat{a}_2 = \text{proj}_S v + \text{proj}_{S^\perp} v$ . Then

$$\begin{aligned} \mathbb{P}[\mathcal{E}, (\mathcal{B}'_{3,12})^c] &= \mathbb{P}[|\langle u, v \rangle| \leq 2\sqrt{d}\alpha \mid \mathcal{E}] \mathbb{P}[\mathcal{E}] \\ (4.10) \quad &\leq \mathbb{P}[|\langle u, \text{proj}_S v \rangle| \leq 2\sqrt{d}\alpha + |\langle u, \text{proj}_{S^\perp} v \rangle| \mid \mathcal{E}] \mathbb{P}[\mathcal{E}]. \end{aligned}$$

When conditioning on  $\mathcal{E}$ ,  $\langle u, \text{proj}_{S^\perp} v \rangle$  is bounded by

$$\begin{aligned} |\langle u, \text{proj}_{S^\perp} v \rangle| &= |u^\top \hat{A}_{>r} \hat{A}_{>r}^\dagger (C_1 \hat{a}_1 - C_2 \hat{a}_2)| \leq \sqrt{dk}\theta \|\hat{A}_{>r}^\dagger (C_1 \hat{a}_1 - C_2 \hat{a}_2)\|_{2/2} \\ (4.11) \quad &\leq \sqrt{dk}\theta\tau M k_l^{-1} m^3. \end{aligned}$$

With (4.11), we can drop the conditioning in (4.10):

$$\begin{aligned}
(4.12) \quad \mathbb{P}[\mathcal{E} \cap (\mathcal{B}'_{3,12})^c] &\leq \mathbb{P}[|\langle u, \text{proj}_S v \rangle| \leq \alpha + \sqrt{dk}\theta\tau M k_l^{-1} m^3] \mathbb{P}[\mathcal{E}] \\
&\leq 2(\alpha + \sqrt{dk}\theta\tau M k_l^{-1} m^3) / (\sqrt{2\pi} \|\text{proj}_S v\|_2) \mathbb{P}[\mathcal{E}] \\
&\leq 2\tau M (\alpha + \sqrt{dk}\theta\tau M k_l^{-1} m^3) \mathbb{P}[\mathcal{E}].
\end{aligned}$$

The last inequality holds because the set  $\{\hat{a}_1, \hat{a}_2, \hat{a}_{r+1}, \dots, \hat{a}_{r+k}\}$  satisfies the robust Kruskal rank condition, and thus

$$\|\text{proj}_S v\|_2 = \|C_1 \hat{a}_1 - C_2 \hat{a}_2 - \hat{A}_{>r} \hat{A}_{>r}^\dagger v\|_2 \geq (\tau M)^{-1} \sqrt{C_1^2 + C_2^2 + \|\hat{A}_{>r}^\dagger v\|_2^2} \geq \sqrt{2} (\tau M)^{-1}.$$

The combination of Lemma B.5 and (4.9) and (4.12) gives the desired probability.  $\blacksquare$

Finally, we are in a place to give the probability that all the events are true for  $x, y$ .

**Lemma 4.10.** *In the setting of Lemma 4.8,  $\mathbb{P}[\mathcal{E}_{1,x}, \mathcal{E}_{1,y}, \mathcal{E}_{2,x}, \mathcal{E}_{2,y}, \mathcal{E}_3] \geq p_1 p_2$ . In particular, the choices  $k_l = \sqrt{2\pi} \tau^{-1} M^{-1} m^3 r^{-1} d^{-1/2} / 64$ ,  $\alpha = \tau^{-1} M^{-1} r^{-2} / 16$ , and  $\theta(r\sqrt{dk}\tau^2 M^2 + 64r^3 \tau^3 M^3 d\sqrt{k}/\sqrt{2\pi}) \leq 1/16$  imply  $\mathbb{P}[\mathcal{E}_{1,x}, \mathcal{E}_{1,y}, \mathcal{E}_{2,x}, \mathcal{E}_{2,y}, \mathcal{E}_3] \geq (\theta\sqrt{d}/8)^{2k} / 256$ .*

*Proof.* The first part follows by combining Lemmas 4.8 and 4.9. For the second part, since  $p_2 \leq p_1$ , the claim follows by using our choices in  $\mathbb{P}[\mathcal{E}_{1,x}, \mathcal{E}_{1,y}, \mathcal{E}_{2,x}, \mathcal{E}_{2,y}, \mathcal{E}_3] \geq p_2^2$ .  $\blacksquare$

At this point, we finished the analysis of the randomness in the first partial tensor decomposition, to recover the first  $r$  components. In the next lemma, we give the probability that random vectors  $x', y'$  satisfy the assumptions of Lemma 4.6. The events will be denoted by  $\mathcal{E}'_{2,x} = \{\forall i \in [k], |\langle x', \hat{a}_{r+i} \rangle| \geq k'_l / m^3\}$ ,  $\mathcal{E}'_{2,y} = \{\forall i \in [k], |\langle y', \hat{a}_{r+i} \rangle| \geq k'_l / m^3\}$ , and  $\mathcal{E}'_3 = \{\forall i \neq j, i, j \in [k], |\langle x', \hat{a}_{r+i} \rangle / \langle y', \hat{a}_{r+i} \rangle - \langle x', \hat{a}_{r+j} \rangle / \langle y', \hat{a}_{r+j} \rangle| \geq \alpha' > 0\}$ .

**Lemma 4.11.** *Let  $x', y'$  be i.i.d. uniformly random in  $\mathcal{S}^{d-1}$ . For  $\hat{a}_{r+1}, \dots, \hat{a}_{r+k}$ , and  $k'_l, \alpha' > 0$ , we have  $\mathbb{P}[\mathcal{E}'_{2,x}, \mathcal{E}'_{2,y}, \mathcal{E}'_3] \geq (1 - k^2 \sqrt{ed} \tau M \alpha' - \sqrt{ed} k k'_l / m^3) (1 - \sqrt{ed} k k'_l / m^3)$ . In particular, the choices  $k'_l = m^3 k^{-1} d^{-1/2} / (4\sqrt{e})$ ,  $\alpha' = \tau^{-1} M^{-1} k^{-2} d^{-1/2} / (4\sqrt{e})$  imply  $\mathbb{P}[\mathcal{E}'_{2,x}, \mathcal{E}'_{2,y}, \mathcal{E}'_3] \geq 3/8$ .*

*Proof.* The first part reuses ideas from the proofs of Lemmas 4.8 and 4.9. We first separate the intersection of events:  $\mathbb{P}[\mathcal{E}'_{2,x} \cap \mathcal{E}'_{2,y} \cap \mathcal{E}'_3] = \mathbb{P}[\mathcal{E}'_{2,x} \cap \mathcal{E}'_3 \mid \mathcal{E}'_{2,y}] \mathbb{P}[\mathcal{E}'_{2,y}] \geq (\mathbb{P}[\mathcal{E}'_3 \mid \mathcal{E}'_{2,y}] - \mathbb{P}[(\mathcal{E}'_{2,x})^c \mid \mathcal{E}'_{2,y}]) \mathbb{P}[\mathcal{E}'_{2,y}]$ . By Lemma B.3,  $\mathbb{P}[\mathcal{E}'_{2,x}]$  and  $\mathbb{P}[\mathcal{E}'_{2,y}]$  are at least  $1 - \sqrt{ed} k k'_l / m^3$ . Also

$$\mathbb{P}[(\mathcal{E}'_3)^c \mid \mathcal{E}'_{2,y}] = \mathbb{E} \left[ \mathbb{P} \left[ \min_{i \neq j, i, j \in [k]} \left| \frac{\langle x', \hat{a}_{r+i} \rangle}{\langle y', \hat{a}_{r+i} \rangle} - \frac{\langle x', \hat{a}_{r+j} \rangle}{\langle y', \hat{a}_{r+j} \rangle} \right| \leq \alpha' \mid y' \right] \mid \mathcal{E}'_{2,y} \right].$$

Consider a uniform upper bound for  $\mathbb{P}[\min_{i \neq j, i, j \in [k]} |\langle x', C'_i \hat{a}_{r+i} - C'_j \hat{a}_{r+j} \rangle| \leq \alpha']$ , where  $|C'_i|, |C'_j|$  are lower bounded by 1. Therefore, again by Lemma B.3, we have  $\mathbb{P}[(\mathcal{E}'_3)^c \mid \mathcal{E}'_{2,y}] \leq k(k-1) \sqrt{ed} \tau M \alpha' / (2\sqrt{2}) \leq k^2 \sqrt{ed} \tau M \alpha'$ . Combining everything gives the desired result. The second part follows directly from our choices of  $k'_l$  and  $\alpha'$ .  $\blacksquare$

**4.4. Putting everything together.** In this subsection, we prove Theorem 3.1.

*Proof of Theorem 3.1.* Without loss of generality, assume  $\pi$  is the identity, and assume for a moment that  $\varepsilon_{in}, \theta$  are small enough so that (1) the assumptions of Lemmas 4.5 and 4.7 are

satisfied; and (2)  $\varepsilon_{4.4}$  and  $\varepsilon_{4.6}$  are smaller than 1 so that we can replace  $\varepsilon_{4.4}^2$  and  $\varepsilon_{4.6}^2$  by  $\varepsilon_{4.4}$  and  $\varepsilon_{4.6}$  in the expression of  $\varepsilon_{4.5}$  and  $\varepsilon_{4.7}$ . We trace the error propagation backwards and show how we can reach  $\varepsilon$  accuracy for the algorithm to terminate while achieving nonnegligible success probability per iteration. The reconstruction error is bounded with [Lemmas 4.4](#) to [4.7](#):

$$\begin{aligned}
(4.13) \quad \|T' - \tilde{T}\|_F &\leq \|\tilde{T} - T\|_F + \sum_{i \in [r+k]} \|a_i^{\otimes 3} - \xi_i \tilde{a}_i^{\otimes 3}\|_F \\
&\leq \varepsilon_{in} + \sum_{i \in [r+k]} \left| \|a_i\|_2^3 - s_i \xi_i \right| \|\tilde{a}_i^{\otimes 3}\|_F + \|\hat{a}_i^{\otimes 3} - s_i^3 \tilde{a}_i^{\otimes 3}\|_F \|a_i\|_2^3 \\
&\leq \varepsilon_{in} + 3rM^3\varepsilon_{4.4} + r\varepsilon_{4.5} + 3kM^3\varepsilon_{4.6} + k\varepsilon_{4.7}.
\end{aligned}$$

Collecting the results from [Lemmas 4.4](#) to [4.7](#), we have

$$(4.14) \quad \begin{aligned}
\varepsilon_{4.4} &= O(\tau^4 M^{10} k r^{5/2} k_l^{-2} \alpha^{-1} (\varepsilon_{in} + \theta)), & \varepsilon_{4.6} &= O(\tau^4 M^7 k^{5/2} r k_l'^{-2} \alpha'^{-1} \varepsilon_{4.5}), \\
\varepsilon_{4.5} &= O(M^3 m^3 r k_l^{-1} \varepsilon_{4.4}), & \varepsilon_{4.7} &= O(M^3 m^3 k k_l'^{-1} \varepsilon_{4.6}).
\end{aligned}$$

With our choices of  $k_l, \alpha, k_l', \alpha'$  in [Lemmas 4.10](#) and [4.11](#), (4.14) can be further written as

$$\begin{aligned}
\varepsilon_{4.4} &= O(\tau^7 M^{13} m^{-6} k r^{13/2} d (\varepsilon_{in} + \theta)), & \varepsilon_{4.6} &= O(\tau^{13} M^{25} m^{-12} k^{11/2} r^{19/2} d^3 (\varepsilon_{in} + \theta)), \\
\varepsilon_{4.5} &= O(\tau^8 M^{17} m^{-6} k r^{17/2} d^{3/2} (\varepsilon_{in} + \theta)), & \varepsilon_{4.7} &= O(\tau^{13} M^{28} m^{-12} k^{13/2} r^{19/2} d^{7/2} (\varepsilon_{in} + \theta)),
\end{aligned}$$

which implies the reconstruction error is bounded by

$$\|T' - \tilde{T}\|_F = O(\tau^{13} M^{28} m^{-12} k^{15/2} r^{19/2} d^{7/2} (\varepsilon_{in} + \theta)).$$

This gives a polynomial  $q(d, r, k, \tau, M, m^{-1}) = \Theta(\tau^{13} M^{28} m^{-12} k^{15/2} r^{19/2} d^{7/2})$ , increasing in every argument, such that if we request that  $\varepsilon_{in} \leq \varepsilon/q(d, r, k, \tau, M, m^{-1})$  and we set  $\theta = \varepsilon/q(d, r, k, \tau, M, m^{-1})$ , then  $\|T' - \tilde{T}\|_F \leq \varepsilon$  (the first termination condition). With this choice, (1) the assumptions of [Lemma 4.10](#) are satisfied; (2) for each iteration, with positive probability the events in [Lemmas 4.10](#) and [4.11](#) happen; and (3) we can take  $\varepsilon_{4.4} = \Theta(\tau^{-6} M^{-15} m^6 r^{-3} d^{-5/2} \varepsilon)$ ,  $\varepsilon_{4.6} = \Theta(M^{-3} k^{-4} d^{-1/2} \varepsilon)$  and they satisfy the assumptions of [Lemmas 4.5](#) and [4.7](#), respectively.

Now we argue that the second termination condition,  $\max_{i \in [r+k]} |\xi_i|^{1/3} \leq 2M$ , holds when the events in [Lemmas 4.10](#) and [4.11](#) happen. Notice that at this point  $|\xi_i|$  is close to  $\|a_i\|_2^3$ . Without loss of generality,  $\max_{i \in [r+k]} |\xi_i|^{1/3} = |\xi_1|^{1/3}$ . Since for all  $x, y > 0$ ,  $|y^{1/3} - x^{1/3}| \leq y^{-2/3}|y - x|$ , we have, for all  $i \in [r+k]$ ,  $|\|a_i\|_2 - |\xi_i|^{1/3}| \leq \|a_i\|_2^{-2} |\|a_i\|_2^3 - |\xi_i||$ , which implies  $|\xi_1|^{1/3} \leq \|a_1\|_2 + |\|a_1\|_2 - |\xi_1|^{1/3}| \leq \|a_1\|_2 + \|a_1\|_2^{-2} \varepsilon \leq M + m \leq 2M$ , where the second inequality comes from  $\varepsilon_{4.5} \leq \varepsilon$  and the third inequality comes from  $\varepsilon \leq \varepsilon_{out} \leq m^3$ . Therefore, the algorithm terminates with a  $2M$ -bounded decomposition with reconstruction error at most  $\varepsilon$ .

Set  $\text{poly}_{3.1}(d, \tau, M) = 2 \text{poly}_{4.2}(2d, \tau, M, 2M, d) \geq 2 \text{poly}_{4.2}(r+k, \tau, M, 2M, d)$ , and set  $\text{poly}'_{3.1}(d, \tau, M, m^{-1}) = q(d, d, d, \tau, M, m^{-1}) \text{poly}_{3.1} \geq q(d, r, k, \tau, M, m^{-1}) \text{poly}_{3.1}$ .<sup>3</sup> When the algorithm terminates, we have

$$\|T - T'\|_F \leq \varepsilon + \varepsilon_{in} \leq \varepsilon + \frac{\varepsilon}{q} \leq \frac{\varepsilon_{out}}{\text{poly}_{3.1}} + \frac{\varepsilon_{out}}{q \text{poly}_{3.1}} \leq \frac{\varepsilon_{out}}{\text{poly}_{4.2}(r+k, \tau, M, 2M, d)}.$$

<sup>3</sup>Recall that  $k \leq r \leq d$  by assumption.

Thus, we can apply [Corollary 4.2](#) and obtain componentwise  $\varepsilon_{out}$  accuracy.

For the running time, in each iteration, steps [3](#) and [7](#) run in time  $\text{poly}(d, \varepsilon^{-1}, \tau, M, m^{-1})$ . Least squares steps [4](#) and [8](#) and the rest take  $\text{poly}(d)$  time. By [Lemmas 4.10](#) and [4.11](#), the success probability per iteration is at least  $3(\theta\sqrt{d}/8)^{2k}/2^{11}$ , which implies that the expected number of iterations is at most  $2^{11}(\theta\sqrt{d}/8)^{-2k}/3$  and the expected running time is at most  $\text{poly}(d^k, \varepsilon^{-k}, \tau^k, M^k, m^{-k})$ . Since  $\varepsilon = \varepsilon_{out}/\text{poly}_{3,1}$ , the expected running time is also at most  $\text{poly}(d^k, 1/\varepsilon_{out}^k, \tau^k, M^k, m^{-k})$ . ■

**5. Blind deconvolution of discrete distribution.** In this section, we provide an application of [Algorithm 2](#): to perform blind deconvolution of an additive mixture model of the form

$$(5.1) \quad X = Z + \eta$$

in  $\mathbb{R}^d$ , where  $Z$  follows a discrete distribution that takes value  $\mu_i$  with probability  $w_i$  for  $i \in [d]$ , and  $\eta$  is an unknown random variable independent of  $Z$  with zero mean, zero 3rd moment, and finite 6th moment.

Our goal is to recover the parameters of  $Z$  when given samples from  $X$ . By estimating the overall mean and translating the samples, we can, without loss of generality, assume that  $\sum_{i \in [d]} w_i \mu_i = 0$  for the rest of this section.

We will see that, under a natural nondegeneracy condition, [Assumption 5.1](#), the parameters of  $Z$  are identifiable from the 3rd cumulant of  $X$ , as the 1st and 3rd moments of  $\eta$  are zero. Let  $K_m(X)$  be the  $m$ th cumulant of  $X$ . By properties of cumulants (see [section 2](#)),

$$(5.2) \quad K_3(X) = K_3(Z) + K_3(\eta) = \sum_{i \in [d]} w_i \mu_i^{\otimes 3}.$$

If the symmetric decomposition of  $K_3(X)$  coincides with [\(5.2\)](#), then the function  $w_i^{1/3} \mu_i$  of the centers  $\mu_i$  and the mixing weights  $w_i$  is identifiable. However, the component vectors satisfy  $\sum_i w_i \mu_i = 0$  (they are *always* linearly dependent), and therefore applying the simultaneous diagonalization algorithm naively has no guarantee.<sup>4</sup> We show that, under the following nondegeneracy condition, our overcomplete tensor decomposition algorithm ([Algorithm 2](#)) works successfully.

*Assumption 5.1.*  $\text{K-rank}_\tau([\mu_1, \dots, \mu_d]) = d - 1$ .

*Remark 5.2.* Note that at this point we are working with a centered mixture ( $\sum_{i \in [d]} w_i \mu_i = 0$ ), and thus the assumption is on the centered mixture. Note also that if  $X = Z + \eta$  is a not necessarily centered mixture, the assumption is satisfied automatically by the centered version of  $X$  when  $Z$  has affinely independent support.

Under [Assumption 5.1](#), we can decompose [\(5.2\)](#) with [Algorithm 2](#). For simplicity, we reformulate the problem: letting  $a_i = \hat{\mu}_i$ , and letting  $\rho_i = \|\mu_i\|_2$ , our goal becomes to decompose  $T = \sum_{i \in [d]} w_i \rho_i^3 a_i^{\otimes 3}$  subject to  $\sum_{i \in [d]} w_i = 1$  and  $\sum_{i \in [d]} w_i \rho_i a_i = 0$ .

We now state our algorithm ([Algorithm 3](#)) for blind deconvolution of discrete distribution.

---

<sup>4</sup>Note that even when the overall mean is nonzero and the means are linearly independent,  $K_3$  still has linearly dependent components, as it is the *central* 3rd moment. If one does not use  $K_3$ , then one loses [\(5.2\)](#).

**Algorithm 3.** Blind deconvolution of discrete distribution.

- Inputs:** i.i.d. samples  $x_1, \dots, x_N$  from mixture  $X$ , error tolerance  $\varepsilon'$ , upper bound  $\rho_{max}$  on  $\|\mu_i\|_2$  for  $i \in [d]$ , lower bound  $w_{min}$  on  $w_i$  for  $i \in [d]$ , robust Kruskal rank threshold  $\tau$ .
- 1: compute the sample 3rd cumulant  $\tilde{T}$  using [Fact A.1](#);
  - 2: invoke [Algorithm 2](#) with error tolerance  $\varepsilon_{5.3} = \varepsilon' / \text{poly}_{5.3}$ , tensor rank  $d$ , and overcompleteness 1 to decompose  $\tilde{T}$  and thus obtain  $\tilde{a}_i \xi_i^{1/3} = \tilde{a}_i$  for  $i \in [d]$ ;
  - 3: set  $\tilde{v}$  to the right singular vector associated with the minimum singular value of  $\tilde{A}$ ;
  - 4: set  $\tilde{w} = \tilde{v}^{3/2} / (\sum_{i \in [d]} \tilde{v}_i^{3/2})$ ,  $\tilde{\mu}_i = \tilde{w}_i^{-1/3} \tilde{a}_i$  for  $i \in [d]$ ;
- Outputs:** estimated mixing weights  $\tilde{w}_1, \dots, \tilde{w}_d$  and estimated means  $\tilde{\mu}_1, \dots, \tilde{\mu}_d$ .

**Theorem 5.3 (correctness of Algorithm 3).** *Let  $X = (X_1, \dots, X_d) = Z + \eta$  be a random vector as in (5.1) satisfying Assumption 5.1. Assume  $0 < w_{min} \leq \min_{i \in [d]} w_i$ ,  $\rho_{max} \geq \max_{i \in [d]} \rho_i$ ,  $0 < \rho_{min} \leq \min_{i \in [d]} \rho_i$ ,  $0 < \varepsilon' \leq \min\{1, w_{min} \rho_{min}^3\}$ , and  $\delta \in (0, 1)$ . There exists a polynomial  $\text{poly}_{5.3}(d, \tau, \rho_{max}, w_{min}^{-1})$  such that if  $\varepsilon_{5.3} = \varepsilon' / \text{poly}_{5.3}$ , then given  $N$  i.i.d. samples of  $X$ , with probability  $1 - \delta$  over the randomness in the samples, [Algorithm 3](#) outputs  $\tilde{\mu}_1, \dots, \tilde{\mu}_d$  and  $\tilde{w}_1, \dots, \tilde{w}_d$  such that for some permutation  $\pi$  of  $[d]$  and for all  $i \in [d]$  we have  $\|\mu_{\pi(i)} - \tilde{\mu}_i\|_2 \leq \varepsilon'$  and  $|w_{\pi(i)} - \tilde{w}_i| \leq \varepsilon'$ . The expected running time over the randomness of [Algorithm 2](#) is at most  $\text{poly}(d, \varepsilon'^{-1}, \delta^{-1}, \tau, \rho_{max}, \rho_{min}^{-1}, w_{min}^{-1}, \max_i \mathbb{E}[X_i^6])$  and will use  $N = \Omega(\varepsilon'^{-2} \delta^{-1} d^{11} \max_{i \in [d]} \mathbb{E}[X_i^6] (\text{poly}'_{3.1}(d, \tau, \rho_{max}, w_{min}^{-1/3} \rho_{min}^{-1}))^2)$  samples.*

The proof of [Theorem 5.3](#) has two parts. First, we show that the 3rd cumulant can be estimated to within  $\varepsilon$  accuracy with polynomially many samples. This follows from a standard argument using  $k$ -statistics. The second part is about the tensor decomposition. Note that [Theorem 3.1](#) guarantees that we can recover  $\tilde{a}_i$  approximately in the direction of  $a_{\pi(i)}$  and  $\xi_i$  close to  $w_{\pi(i)} \rho_{\pi(i)}^3$  for some permutation  $\pi$ . However, we are not finished yet, as our goal is to recover both the centers and the mixing weights. Therefore we need to decouple  $w_i$  and  $\rho_i$  from  $w_i \rho_i^3$ , which corresponds to steps 3 and 4 in [Algorithm 3](#).

**3rd cumulant estimation.** The details are in [Appendix A](#); we only give the main result here.

**Lemma 5.4 (estimation of the 3rd cumulant).** *Let  $T, \tilde{T}$  be the 3rd cumulant of  $X = (X_1, \dots, X_d)$  and its unbiased estimate ( $k$ -statistic) using [Fact A.1](#), respectively. Given any  $\varepsilon, \delta \in (0, 1)$ , and  $N = \Omega(d^9 \varepsilon^{-2} \delta^{-1} \max_{i \in [d]} \mathbb{E}[X_i^6])$ , with probability  $1 - \delta$  we have  $\|T - \tilde{T}\|_F \leq \varepsilon$ .*

*Proof.* Apply [Lemma A.3](#) with accuracy  $\varepsilon/d^3$  and failure probability  $\delta/d^3$ , and take the union bound over  $d^3$  entries to see that  $N = \Omega(d^9 \varepsilon^{-2} \delta^{-1} \max_{i \in [d]} \mathbb{E}[X_i^6])$  samples are sufficient. ■

**Decoupling.** We will decouple the mixing weights  $w_i$  and the norms  $\rho_i$  after we decompose the tensor  $\tilde{T}$ . As  $\mathbb{E}[X] = 0$ , the true parameters satisfy  $\sum_{i \in [d]} w_i \rho_i a_i = 0$ , which can be reformulated as a linear system

$$(5.3) \quad AD_{w_i \rho_i^3}^{1/3} w^{2/3} = 0,$$

where  $D_{w_i \rho_i^3} = \text{diag}(w_i \rho_i^3)$  and  $A$  contains  $a_i$ s as columns. To decouple these parameters in the noiseless setting, one only needs to solve this system under the constraint that  $w$  is a

probability vector. As  $\text{rank}(A) = d - 1$ ,  $w$  will be uniquely determined. In other words,  $w^{2/3}$  lies in the direction of the right singular vector associated with the only zero singular value. It is natural then to recover the weights using our approximations to terms in the linear system, namely in the direction of the right singular vector associated with the minimum singular value of  $\tilde{A}D_\xi^{1/3}$ , where  $\tilde{A} = [\tilde{a}_1, \dots, \tilde{a}_d]$  and  $D_\xi = \text{diag}(\xi_i)$ . The following theorem guarantees this will work.

**Theorem 5.5 (decoupling).** *Let  $0 < w_{\min} \leq \min_{i \in [d]} w_i$ , and let  $\rho_{\max} \geq \max_{i \in [d]} \|\mu_i\|_2$ . Suppose the outputs of step 2 in Algorithm 3, namely  $\xi_1, \dots, \xi_d$  and  $\tilde{A} = [\tilde{a}_1, \dots, \tilde{a}_d]$ , satisfy Theorem 3.1 with  $\varepsilon_{\text{out}} < w_{\min}^{4/3}/(24d\tau)$  and permutation  $\pi$ . One can choose positive right singular vectors  $v, \tilde{v}$  associated with the minimum singular values of  $AD_{w_i\rho_i^3}^{1/3}, \tilde{A}D_\xi^{1/3}$ , respectively. Define  $\tilde{w} = \tilde{v}^{3/2}/\sum_{i \in [d]} \tilde{v}_i^{3/2}$  and  $\tilde{\rho}_i = (\xi_i/\tilde{w}_i)^{1/3}$ . Then  $|w_{\pi(i)} - \tilde{w}_i| \leq 12w_{\min}^{-1/3}d\tau\varepsilon_{\text{out}}$  and  $|\rho_{\pi(i)} - \tilde{\rho}_i| \leq 48w_{\min}^{-4/3}\rho_{\max}d\tau\varepsilon_{\text{out}}$ .*

*Proof.* We start by showing that  $v, \tilde{v}$ , and  $\tilde{w}$  are well-defined. Since  $w^{2/3}$  is a solution to (5.3) and  $AD_{w_i\rho_i^3}^{1/3}$  is of rank  $d - 1$ , we pick  $v = w^{2/3}/\|w^{2/3}\|_2$ . To show that  $\tilde{v}$  is well-defined, first we bound the singular values and vectors of  $\tilde{A}D_\xi^{1/3}$ . Let  $\tilde{\sigma}_i = \sigma_i(\tilde{A}D_\xi^{1/3})$ . By Theorem B.1,

$$(5.4) \quad \tilde{\sigma}_d \leq \|AD_{w_i\rho_i^3}^{1/3} - \tilde{A}D_\xi^{1/3}\|_2 \leq \sqrt{d}\varepsilon_{\text{out}} < w_{\min}^{4/3}/(24\sqrt{d}\tau).$$

To obtain the deviation in the singular vectors, we first show that  $\tilde{\sigma}_1, \dots, \tilde{\sigma}_{d-1}$  are bounded away from zero. Let  $\Sigma_1 = \text{diag}(\sigma_1(AD_{w_i\rho_i^3}^{1/3}), \dots, \sigma_{d-1}(AD_{w_i\rho_i^3}^{1/3}))$ ,  $\tilde{\Sigma}_1 = \text{diag}(\tilde{\sigma}_1, \dots, \tilde{\sigma}_{d-1})$ , and  $\Delta = w_{\min}^{1/3}/(2\tau)$ . Suppose  $\hat{\sigma}_{d-1}$  is the least singular value of the matrix obtained by deleting the first column of  $AD_{w_i\rho_i^3}^{1/3}$ ; then it follows that  $\sigma_{d-1}(AD_{w_i\rho_i^3}^{1/3}) \geq \hat{\sigma}_{d-1} \geq w_{\min}^{1/3}/\tau$ , where the first inequality follows from the interlacing property of singular values of a matrix and its submatrix obtained by deleting any column, and the second inequality comes from Assumption 5.1. The minimum diagonal term in  $\tilde{\Sigma}_1$  satisfies

$$\min_i (\tilde{\Sigma}_1)_{ii} \geq \sigma_{d-1}(AD_{w_i\rho_i^3}^{1/3}) - \sqrt{d}\varepsilon_{\text{out}} \geq \frac{w_{\min}^{1/3}}{\tau} - \frac{w_{\min}^{4/3}}{24\sqrt{d}\tau} \geq \frac{w_{\min}^{1/3}}{2\tau} = \Delta.$$

Therefore by Theorem B.2 with  $\Sigma_2 = 0$ , we have for the singular vectors,<sup>5</sup>

$$\|v - \tilde{v}\|_2 \leq \sqrt{2d}\varepsilon_{\text{out}}/\Delta = 2\sqrt{2d}w_{\min}^{-1/3}\tau\varepsilon_{\text{out}}.$$

We get  $\tilde{v}_i \geq v_i - 2\sqrt{2d}w_{\min}^{-1/3}\tau\varepsilon_{\text{out}} \geq w_{\min}^{2/3} - 2\sqrt{2d}w_{\min}^{-1/3}\tau\varepsilon_{\text{out}} > 0$ , where the second inequality follows from  $\sum_{i \in [d]} v_i^{3/2} \geq \sum_{i \in [d]} v_i^2 = 1$ . Hence  $\tilde{v}$  also has positive entries and  $\tilde{w}$  is well-defined.

<sup>5</sup>Note that even though Theorem B.2 gives the angle between the subspaces spanned by the first  $d - 1$  right singular vectors of  $AD_{w_i\rho_i^3}^{1/3}$  and their perturbed counterparts, the same bound applies to the orthogonal complement, spanned by  $v$ .



We now derive the bounds on the mixing weights and norms. Without loss of generality,  $\pi$  is the identity. The mixing weight error is bounded by

$$(5.5) \quad \begin{aligned} \|\tilde{w} - w\|_2 &= \left\| \frac{\tilde{v}^{3/2}}{\sum_{i \in [d]} \tilde{v}_i^{3/2}} - \frac{v^{3/2}}{\sum_{i \in [d]} v_i^{3/2}} \right\|_2 \\ &\leq \frac{\|\tilde{v}^{3/2} - v^{3/2}\|_2}{\sum_{i \in [d]} v_i^{3/2}} + \frac{\|\tilde{v}^{3/2}\|_2}{(\sum_{i \in [d]} v_i^{3/2})(\sum_{i \in [d]} \tilde{v}_i^{3/2})} \left| \sum_{i \in [d]} (v_i^{3/2} - \tilde{v}_i^{3/2}) \right|. \end{aligned}$$

We bound each term in (5.5) below since  $\tilde{v}, v$  both have entries in  $(0, 1]$ :  $\sum_{i \in [d]} v_i^{3/2} \geq \|v\|_2^2 = 1$ ,  $\sum_{i \in [d]} \tilde{v}_i^{3/2} \geq \|\tilde{v}\|_2^2 = 1$ , and  $\|\tilde{v}^{3/2}\|_2 = (\sum_{i \in [d]} \tilde{v}_i^3)^{1/2} \leq \|\tilde{v}\|_2 = 1$ . Moreover,

$$\|\tilde{v}^{3/2} - v^{3/2}\|_2 = \left( \sum_{i \in [d]} (\tilde{v}_i^{3/2} - v_i^{3/2})^2 \right)^{1/2} \leq \frac{3}{2} \left( \sum_{i \in [d]} (\tilde{v}_i - v_i)^2 \right)^{1/2} = \frac{3}{2} \|\tilde{v} - v\|_2,$$

$$\left| \sum_{i \in [d]} (v_i^{3/2} - \tilde{v}_i^{3/2}) \right| \leq \sum_{i \in [d]} |v_i^{3/2} - \tilde{v}_i^{3/2}| \leq \frac{3}{2} \sum_{i \in [d]} |v_i - \tilde{v}_i| \leq \frac{3\sqrt{d}}{2} \|\tilde{v} - v\|_2,$$

where the above two inequalities follow from  $|x^{3/2} - y^{3/2}| \leq 3|x - y|/2$  for  $x, y \in [0, 1]$ .

We obtain the following bound on the error in mixing weights:

$$(5.6) \quad \|\tilde{w} - w\|_2 \leq (3/2)(1 + \sqrt{d})\|\tilde{v} - v\|_2 \leq 3\sqrt{2}w_{\min}^{-1/3}(d + \sqrt{d})\tau\varepsilon_{out} \leq 12w_{\min}^{-1/3}d\tau\varepsilon_{out}.$$

Notice that our assumption on  $\varepsilon_{out}$  guarantees that  $\tilde{w}_i \geq w_{\min}/2$ , and therefore the error in the norm is bounded by

$$\begin{aligned} |\tilde{\rho}_i - \rho_i| &= |(\xi_i/\tilde{w}_i)^{1/3} - \rho_i| \leq \tilde{w}_i^{-1/3} (|\xi_i^{1/3} - w_i^{1/3}\rho_i| + \rho_i|w_i^{1/3} - \tilde{w}_i^{1/3}|) \\ &\leq \tilde{w}_i^{-1/3} (\varepsilon_{out} + \rho_{\max}|w_i^{1/3} - \tilde{w}_i^{1/3}|) \leq (2w_{\min}^{-1})^{1/3} (\varepsilon_{out} + \rho_{\max}w_{\min}^{-2/3}|w_i - \tilde{w}_i|) \\ &\leq 48w_{\min}^{-4/3} \rho_{\max}d\tau\varepsilon_{out}, \end{aligned}$$

where the second inequality comes from [Theorem 3.1](#), the third inequality comes from the fact that  $|x^{1/3} - y^{1/3}|/|x - y| \leq y^{-2/3}$  for any  $x, y > 0$ , and the last one follows from (5.6). ■

We are now ready to prove [Theorem 5.3](#).

*Proof of Theorem 5.3.* Set the arguments of  $\text{poly}'_{3.1}, \text{poly}_{3.1}$  to  $(d, \tau, \rho_{\max}, w_{\min}^{-1/3} \rho_{\min}^{-1})$  and  $(d, \tau, \rho_{\max})$ , respectively. Assume for a moment that  $N$  is large enough so that  $T$  in step 1 satisfies  $\|T - \tilde{T}\|_F \leq \varepsilon_{out}/(\text{poly}'_{3.1})$  and we can apply [Theorem 3.1](#). We start by verifying that we can apply [Theorem 5.5](#). Set  $\text{poly}_{5.3} = 49w_{\min}^{-4/3} \max\{\rho_{\max}, 1\}d\tau \text{poly}_{3.1}$ . By [Theorem 3.1](#), our choice of  $\varepsilon_{5.3}$  guarantees that the output error of step 2 in [Algorithm 3](#) is  $\varepsilon_{out} = \varepsilon_{5.3} \text{poly}_{3.1} = \varepsilon'/(49w_{\min}^{-4/3} \max\{\rho_{\max}, 1\}d\tau) < w_{\min}^{4/3}/(24d\tau)$  (using our assumption  $\varepsilon' \leq 1$ ). We now bound our estimation error for  $\|\mu_i\|_2$  and  $w_i$  with [Theorem 5.5](#). Assuming the permutation is the identity, we have for  $i \in [d]$ ,

$$\begin{aligned} \|\mu_i - \tilde{\mu}_i\|_2 &\leq |\rho_i - \tilde{\rho}_i| \|\tilde{a}_i\|_2 + \rho_i \|a_i - \tilde{a}_i\|_2 \leq (48w_{\min}^{-4/3} \rho_{\max}d\tau + \rho_{\max})\varepsilon_{out} \leq \varepsilon', \\ |w_i - \tilde{w}_i| &\leq 12w_{\min}^{-1/3}d\tau\varepsilon_{out} \leq \varepsilon'. \end{aligned}$$

Next, we derive the sample complexity: we need

$$\|T - \tilde{T}\|_F \leq \varepsilon_{in} \leq \varepsilon_{out}/(\text{poly}'_{3,1}) = \varepsilon' \text{poly}_{3,1}/(\text{poly}'_{3,1} \text{poly}_{5,3}).$$

By [Lemma 5.4](#),  $N = \Omega(\varepsilon'^{-2} \delta^{-1} d^{11} \max_{i \in [d]} \mathbb{E}[X_i^6] (\text{poly}'_{3,1})^2)$  many samples are sufficient for  $\varepsilon_{in}$  to meet the assumption. Since  $N$  is polynomial in  $\delta^{-1}$  and  $\max_{i \in [d]} \mathbb{E}[X_i^6]$ , the expected running time will also be polynomial in them.  $\blacksquare$

**6. Parameter estimation of Gaussian mixture models (GMM).** In this section, we consider a specific family of mixture models, namely GMM with identical but unknown covariance matrices. The model is as in [\(5.1\)](#), where  $\eta \sim \mathcal{N}(0, \Sigma)$ . Our goal is to approximate all parameters of the mixture:  $\Sigma$ ,  $w_i$ s, and  $\mu_i$ s. Again, suppose [Assumption 5.1](#) holds and the mean of the mixture is zero (by translating the samples as in [section 5](#)). [Algorithm 3](#) guarantees that we can recover the mixing weights  $w_i$ s and centers  $\mu_i$ s of  $Z$ . To recover  $\Sigma$ , notice that since the mean is zero,  $\text{cov}(X) = \mathbb{E}[XX^\top] = \sum_{i \in [d]} w_i \mu_i \mu_i^\top + \Sigma$ . The covariance matrix can be approximated then by taking the difference between the sample second moment of  $X$  and the second moment of the reconstructed discrete distribution. We make this precise in [Algorithm 4](#) and [Theorem 6.1](#).

---

**Algorithm 4.** Parameter estimation for GMM.

---

- Inputs:** i.i.d. samples  $x_1, \dots, x_N$  from mixture  $X$ , error tolerance  $\varepsilon''$ , upper bound  $\rho_{max}$  on  $\|\mu_i\|_2$  for  $i \in [d]$ , lower bound  $w_{min}$  on  $w_i$  for  $i \in [d]$ , robust Kruskal rank threshold  $\tau$ .
- 1: invoke [Algorithm 3](#) with samples from  $X$  and parameters  $\varepsilon' = \varepsilon''/\text{poly}_{6,1}$ ,  $\rho_{max}$ ,  $w_{min}$ ,  $\tau$  to get  $\tilde{w}_i$  and  $\tilde{\mu}_i$  for  $i \in [d]$ ;
  - 2: set  $\tilde{\Sigma} = \frac{1}{N} \sum_{j \in [N]} x_j x_j^\top - \sum_{i \in [d]} \tilde{w}_i \tilde{\mu}_i \tilde{\mu}_i^\top$ ;
- Outputs:** estimated covariance matrix  $\tilde{\Sigma}$ , mixing weights and means  $\tilde{w}_i, \tilde{\mu}_i : i \in [d]$ .
- 

**Theorem 6.1 (correctness of Algorithm 4).** *Let  $X$  be a GMM with identical but unknown covariance matrices satisfying [Assumption 5.1](#). Assume  $0 < w_{min} \leq \min_{i \in [d]} w_i$ ,  $\rho_{max} \geq \max_{i \in [d]} \rho_i$ ,  $0 < \rho_{min} \leq \min_{i \in [d]} \rho_i$ ,  $0 < \varepsilon'' \leq \min\{1, w_{min} \rho_{min}^3\}$ , and  $\delta \in (0, 1)$ . There exists a polynomial  $\text{poly}_{6,1}(d, \rho_{max})$  such that if  $\varepsilon' = \varepsilon''/\text{poly}_{6,1}$ , then given  $N$  i.i.d. samples of  $X$  and with probability  $1 - \delta$  over the randomness in the samples, [Algorithm 4](#) outputs  $\tilde{\mu}_1, \dots, \tilde{\mu}_d$ ,  $\tilde{w}_1, \dots, \tilde{w}_d$ , and  $\tilde{\Sigma}$  such that for some permutation  $\pi$  of  $[d]$  and for all  $i \in [d]$ :  $\|\tilde{\Sigma} - \Sigma\|_F \leq \varepsilon''$ ,  $|w_{\pi(i)} - \tilde{w}_i| \leq \varepsilon'$ , and  $\|\mu_{\pi(i)} - \tilde{\mu}_i\|_2 \leq \varepsilon'$ . The expected running time over the randomness of [Algorithm 2](#) is at most  $\text{poly}(d, \varepsilon''^{-1}, \delta^{-1}, \tau, \rho_{max}, \rho_{min}^{-1}, w_{min}^{-1}, \max_{i \in [d]} \Sigma_{ii}^3)$  and will use  $N = \Omega(\varepsilon''^{-2} \delta^{-1} d^{13} \max_{i \in [d]} \Sigma_{ii}^3 (\text{poly}'_{3,1}(d, \tau, \rho_{max}, w_{min}^{-1/3}, \rho_{min}^{-1}))^2)$  samples.*

*Proof.* Let  $\text{poly}_{6,1}(d, \rho_{max}) = 1 + d\rho_{max}^2 + 2d(2\rho_{max} + 1)$ . By [Theorem 5.3](#), with probability  $1 - \delta$ , [Algorithm 3](#) will output the estimated mixing weights  $\tilde{w}_i$  and means  $\tilde{\mu}_i$  within  $\varepsilon'$  additive accuracy. The sample complexity and running time follow therein, where we have  $\max_{i \in [d]} \mathbb{E}[X_i^6] = \max_{i \in [d]} 15\Sigma_{ii}^3$  for GMM.

Next, we bound the error in the covariance matrix. Note that when the number of samples guarantees that  $K_3(X)$  is estimated to  $\varepsilon_{in}$  accuracy with probability  $1 - \delta$ , it can also guarantee  $\text{cov}(X)$  is estimated to  $\varepsilon_{in}$  accuracy with probability  $1 - \delta$  since the latter takes

$\Omega(d^6 \varepsilon_{in}^{-2} \delta^{-1} \max_{i \in [d]} \Sigma_{ii}^2)$  many samples by an argument similar to [Lemmas A.2](#) and [A.3](#). So

$$\begin{aligned} \|\tilde{\Sigma} - \Sigma\|_F &= \left\| \frac{1}{N} \sum_{j \in [N]} x_j x_j^\top - \sum_{i \in [d]} \tilde{w}_i \tilde{\mu}_i \tilde{\mu}_i^\top - \Sigma \right\|_F \\ &\leq \left\| \frac{1}{N} \sum_{j \in [N]} x_j x_j^\top - \text{cov}(X) \right\|_F + \sum_{i \in [d]} |w_i - \tilde{w}_i| \|\mu_i \mu_i^\top\|_F + \sum_{i \in [d]} \tilde{w}_i \|\mu_i \mu_i^\top - \tilde{\mu}_i \tilde{\mu}_i^\top\|_F \\ &\leq \varepsilon_{in} + d \rho_{max}^2 \varepsilon' + \sum_{i \in [d]} (w_i + \varepsilon') (2 \|\mu_i\|_2 + \varepsilon') \varepsilon' \leq \text{poly}_{6.1} \varepsilon' \leq \varepsilon'', \end{aligned}$$

where the second to last inequality follows from bounding  $\varepsilon_{in}$  by  $\varepsilon'$  and  $w_i, \varepsilon'$  by 1.  $\blacksquare$

**Appendix A. Estimating cumulants.** In this section, we provide technical details about the unbiased estimators of cumulants, called  $k$ -statistics. They are the unbiased estimator for cumulants with the minimum variance and are long studied in the statistics community. We provide the formula for the 3rd  $k$ -statistic given in [\[32, Chapter 4\]](#) here.

**Fact A.1.** *Given i.i.d. samples  $x_1, \dots, x_N$  of random vector  $X$ , the  $k$ -statistic for the 3rd cumulant of  $X$  is  $k_3(r, s, t) = \frac{1}{N} \sum_{i, j, k \in [N]} \phi^{(ijk)}(x_i)_r (x_j)_s (x_k)_t$ , where  $r, s, t$  are the position indices in the tensor, and  $\phi^{(ijk)}$  is a family of coefficients defined in the following way: it is invariant under permutation of indices, and for distinct  $i, j, k \in [N]$ ,*

$$(A.1) \quad \phi^{(iii)} = \frac{1}{N}, \quad \phi^{(iij)} = -\frac{1}{N-1}, \quad \phi^{(ijk)} = \frac{2}{(N-1)(N-2)}.$$

To obtain the entrywise concentration bound for  $k_3$ , we begin by bounding the variance of each entry in  $k_3$ .

**Lemma A.2.** *Let  $X$  follow a distribution as in [\(5.1\)](#). The 3rd  $k$ -statistics  $k_3$  of  $X$  satisfies  $\text{Var}(k_3(r, s, t)) = O(\max_{t \in [d]} \mathbb{E}[X_t^6]/N)$ .*

*Proof.* An essentially identical result for the 4th cumulant is shown in [\[4, Lemma 4\]](#). The argument here is the same. We provide a proof in the supplementary materials ([supplementary.pdf \[local/web 206KB\]](#)).  $\blacksquare$

Using Chebyshev's inequality yields the following sample bound immediately.

**Lemma A.3.** *Given  $\varepsilon, \delta \in (0, 1)$ , the entrywise error between  $k_3$  and  $K_3(X)$  is at most  $\varepsilon$  with probability at least  $1 - \delta$  when using  $N \geq \Omega(\varepsilon^{-2} \delta^{-1} \max_{t \in [d]} \mathbb{E}[X_t^6])$  samples.*

## Appendix B. Technical lemmas.

**B.1. Perturbed SVD bounds.** We state Wedin's theorem, a "sin( $\theta$ ) theorem" for perturbed singular vectors as well as Weyl's inequality for SVD. The following results are from [\[35, 36\]](#).

**Theorem B.1 (Weyl's inequality).** *Let  $A, E \in \mathbb{R}^{d_1 \times d_2}$  with  $d_1 \geq d_2$ . Denote the singular values in nonincreasing order of  $A$  and  $A+E$  by  $\sigma_i$  and  $\tilde{\sigma}_i$ , respectively. Then  $|\sigma_i - \tilde{\sigma}_i| \leq \|E\|_2$ .*

**Theorem B.2 (Wedin).** *With the notation from Theorem B.1, let a singular value decomposition of  $A$  be*

$$[U_1, U_2, U_3]^\top A [V_1, V_2] = \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \\ 0 & 0 \end{bmatrix},$$

where the singular values can be in arbitrary order. Let the perturbed version be

$$[\tilde{U}_1, \tilde{U}_2, \tilde{U}_3]^\top (A + E) [\tilde{V}_1, \tilde{V}_2] = \begin{bmatrix} \tilde{\Sigma}_1 & 0 \\ 0 & \tilde{\Sigma}_2 \\ 0 & 0 \end{bmatrix}.$$

Let  $\Phi$  be the matrix of canonical angles between the column spaces of  $U_1$  and  $\tilde{U}_1$ , and let  $\Theta$  be that of  $V_1$  and  $\tilde{V}_1$ , respectively. Let  $\delta = \min\{\min_i \tilde{\Sigma}_{1,ii}, \min_{i,j} |\tilde{\Sigma}_{1,ii} - \Sigma_{2,jj}|\}$ . Then  $\sqrt{\|\sin \Phi\|_2^2 + \|\sin \Theta\|_2^2} \leq \sqrt{2} \|E\|_2 / \delta$ .

## B.2. Probability tail bounds.

**Lemma B.3** (see [9, 20]). *Suppose  $\delta \in (0, 1)$ ,  $M \in \mathbb{R}^{d \times d}$ ,  $Q$  is a finite subset of  $\mathbb{R}^d$ , and  $X$  is a uniformly random vector in  $\mathcal{S}^{d-1}$ . Then  $\mathbb{P}[\min_{q \in Q} |\langle X, Mq \rangle| \geq \frac{\delta \min_{q \in Q} \|Mq\|_2}{\sqrt{ed|Q|}}] \geq 1 - \delta$ .*

For the next lemma, we need the Gaussian correlation inequality.

**Theorem B.4** (Gaussian correlation inequality [26, 33]). *For any convex centrally symmetric sets  $K, L$  in  $\mathbb{R}^d$  and any centered Gaussian measure  $\mu$  on  $\mathbb{R}^d$ , we have  $\mu(K \cap L) \geq \mu(K)\mu(L)$ .*

**Lemma B.5** (see [23, 34]). *Let  $X \in \mathbb{R}^d$  be a standard Gaussian random vector, let  $a_1, \dots, a_k \in \mathcal{S}^{d-1}$ , and let  $t \in [0, 1]$ . Then  $\mathbb{P}[(\forall i) |\langle X, a_i \rangle| \leq t] \geq (t/4)^k$ .*

*Proof.* The claim follows immediately from Theorem B.4 and the fact that the one-dimensional standard Gaussian density in  $[-1, 1]$  is at least  $(2\pi e)^{-1/2} \geq 1/8$ . ■

**Acknowledgments.** We would like to thank Nina Amenta, Jesús De Loera, Shuyang Ling, Naoki Saito, and James Sharpnack for helpful discussions.

## REFERENCES

- [1] A. ANANDKUMAR, R. GE, AND M. JANZAMIN, *Guaranteed Non-orthogonal Tensor Decomposition via Alternating Rank-1 Updates*, preprint, <https://arxiv.org/abs/1402.5180>, 2014.
- [2] A. ANANDKUMAR, R. GE, AND M. JANZAMIN, *Learning overcomplete latent variable models through tensor methods*, in Proceedings of the 28th Conference on Learning Theory (COLT 2015), Paris, France, 2015, JMLR Workshop Conf. Proc. 40, 2015, pp. 36–112, <http://proceedings.mlr.press/v40/Anandkumar15.html>.
- [3] J. ANDERSON, M. BELKIN, N. GOYAL, L. RADEMACHER, AND J. VOSS, *The more, the merrier: The blessing of dimensionality for learning large Gaussian mixtures*, in Proceedings of the 27th Conference on Learning Theory (COLT 2014), pp. 1135–1164.
- [4] M. BELKIN, L. RADEMACHER, AND J. R. VOSS, *Blind signal separation in the presence of Gaussian noise*, in Proceedings of the 26th Annual Conference on Learning Theory (COLT 2013), Princeton University, Princeton, NJ, USA, 2013, JMLR Workshop Conf. Proc. 30, S. Shalev-Shwartz and I. Steinwart, eds., 2013, pp. 270–287, <http://proceedings.mlr.press/v30/Belkin13.html>.

- [5] A. BHASKARA, M. CHARIKAR, A. MOITRA, AND A. VIJAYARAGHAVAN, *Smoothed analysis of tensor decompositions*, in Proceedings of the Forty-Sixth Annual ACM Symposium on Theory of Computing (STOC '14), New York, NY, USA, 2014, pp. 594–603, <https://doi.org/10.1145/2591796.2591881>.
- [6] A. BHASKARA, M. CHARIKAR, AND A. VIJAYARAGHAVAN, *Uniqueness of tensor decompositions with applications to polynomial identifiability*, in Proceedings of the Conference on Learning Theory (COLT 2014), pp. 742–778.
- [7] J. CARDOSO, *Super-symmetric decomposition of the fourth-order cumulant tensor. Blind identification of more sources than sensors*, in Proceedings of the 1991 IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP '91), Toronto, ON, Canada, 1991, pp. 3109–3112, <https://doi.org/10.1109/ICASSP.1991.150113>.
- [8] P. COMON AND C. JUTTEN, *Handbook of Blind Source Separation: Independent Component Analysis and Applications*, 1st ed., Academic Press, New York, 2010.
- [9] S. DASGUPTA AND A. GUPTA, *An elementary proof of a theorem of Johnson and Lindenstrauss*, Random Structures Algorithms, 22 (2003), pp. 60–65.
- [10] I. DIAKONIKOLAS, S. B. HOPKINS, D. KANE, AND S. KARMALKAR, *Robustly Learning Any Clusterable Mixture of Gaussians*, preprint, <https://arxiv.org/abs/2005.06417>, 2020.
- [11] I. DOMANOV AND L. DE LATHAUWER, *Canonical polyadic decomposition of third-order tensors: Reduction to generalized eigenvalue decomposition*, SIAM J. Matrix Anal. Appl., 35 (2014), pp. 636–660, <https://doi.org/10.1137/130916084>.
- [12] I. DOMANOV AND L. DE LATHAUWER, *Canonical polyadic decomposition of third-order tensors: Relaxed uniqueness conditions and algebraic algorithm*, Linear Algebra Appl., 513 (2017), pp. 342–375, <https://doi.org/10.1016/j.laa.2016.10.019>.
- [13] R. GE, Q. HUANG, AND S. M. KAKADE, *Learning mixtures of Gaussians in high dimensions*, in Proceedings of the Forty-Seventh Annual ACM Symposium on Theory of Computing, 2015, pp. 761–770.
- [14] R. GE AND T. MA, *Decomposing overcomplete 3rd order tensors using sum-of-squares algorithms*, in Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques, LIPIcs. Leibniz Int. Proc. Inform. 40, N. Garg, K. Jansen, A. Rao, and J. D. P. Rolim, eds., Schloss Dagstuhl. Leibniz-Zent. Inform., Wadern, Germany, 2015, pp. 829–849, <https://doi.org/10.4230/LIPIcs.APPROX-RANDOM.2015.829>.
- [15] R. GE AND T. MA, *On the optimization landscape of tensor decompositions*, in Advances in Neural Information Processing Systems: Annual Conference on Neural Information Processing Systems, Vol. 30, Long Beach, CA, USA, I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, eds., MIT Press, Cambridge, MA, 2017, pp. 3653–3663, <http://papers.nips.cc/paper/6956-on-the-optimization-landscape-of-tensor-decompositions>.
- [16] N. GOYAL, S. VEMPALA, AND Y. XIAO, *Fourier PCA and robust tensor decomposition*, in Proceedings of the Forty-Sixth Annual ACM Symposium on Theory of Computing, 2014, pp. 584–593.
- [17] N. GOYAL, S. S. VEMPALA, AND Y. XIAO, *Fourier PCA*, preprint, <https://arxiv.org/abs/1306.5825v4>, 2013.
- [18] S. B. HOPKINS, T. SCHRAMM, AND J. SHI, *A robust spectral algorithm for overcomplete tensor decomposition*, in Proceedings of the Conference on Learning Theory (COLT 2019), Phoenix, AZ, USA, Proc. Mach. Learn. Res. 99, A. Beygelzimer and D. Hsu, eds., 2019, pp. 1683–1722, <http://proceedings.mlr.press/v99/hopkins19b.html>.
- [19] S. B. HOPKINS, T. SCHRAMM, J. SHI, AND D. STEURER, *Fast spectral algorithms from sum-of-squares proofs: Tensor decomposition and planted sparse vectors*, in Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing (STOC 2016), Cambridge, MA, USA, 2016, D. Wichs and Y. Mansour, eds., 2016, pp. 178–191, <https://doi.org/10.1145/2897518.2897529>.
- [20] D. HSU AND S. M. KAKADE, *Learning mixtures of spherical Gaussians: Moment methods and spectral decompositions*, in Proceedings of the 4th ACM Conference on Innovations in Theoretical Computer Science, 2013, pp. 11–20.
- [21] M. JANZAMIN, R. GE, J. KOSSAIFI, AND A. ANANDKUMAR, *Spectral learning on matrices and tensors*, Found. Trends Mach. Learn., 12 (2019), pp. 393–536, <https://doi.org/10.1561/22000000057>.
- [22] H. JIA AND S. VEMPALA, *Robustly Clustering a Mixture of Gaussians*, preprint, <https://arxiv.org/abs/1911.11838>, 2019.
- [23] C. G. KHATRI, *On certain inequalities for normal distributions and their applications to simultaneous con-*

- fidence bounds*, Ann. Math. Statist., 38 (1967), pp. 1853–1867, <http://www.jstor.org/stable/2238663>.
- [24] T. G. KOLDA, *Will the real Jennrich’s Algorithm please stand up?*, <https://www.mathsci.ai/post/jennrich/>, accessed 2021–12–21.
- [25] J. B. KRUSKAL, *Three-way arrays: Rank and uniqueness of trilinear decompositions, with application to arithmetic complexity and statistics*, Linear Algebra Appl., 18 (1977), pp. 95–138, [https://doi.org/10.1016/0024-3795\(77\)90069-6](https://doi.org/10.1016/0024-3795(77)90069-6).
- [26] R. LATAŁA AND D. MATLAK, *Royen’s proof of the Gaussian correlation inequality*, in Geometric Aspects of Functional Analysis, Springer, Cham, 2017, pp. 265–275.
- [27] L. D. LATHAUWER, J. CASTAING, AND J. CARDOSO, *Fourth-order cumulant-based blind identification of underdetermined mixtures*, IEEE Trans. Signal Process., 55 (2007), pp. 2965–2973, <https://doi.org/10.1109/TSP.2007.893943>.
- [28] S. E. LEURGANS, R. T. ROSS, AND R. B. ABEL, *A decomposition for three-way arrays*, SIAM J. Matrix Anal. Appl., 14 (1993), pp. 1064–1083, <https://doi.org/10.1137/0614071>.
- [29] A. LEVIN, Y. WEISS, F. DURAND, AND W. T. FREEMAN, *Understanding blind deconvolution algorithms*, IEEE Trans. Pattern Anal. Mach. Intell., 33 (2011), pp. 2354–2367.
- [30] T. MA, J. SHI, AND D. STEURER, *Polynomial-time tensor decompositions with sum-of-squares*, in Proceedings of the 57th Annual IEEE Symposium on Foundations of Computer Science (FOCS 2016), 2016, New Brunswick, NJ, USA, I. Dinur, ed., 2016, pp. 438–446, <https://doi.org/10.1109/FOCS.2016.54>.
- [31] T. MA, J. SHI, AND D. STEURER, *Polynomial-time tensor decompositions with sum-of-squares*, in Proceedings of the 57th Annual IEEE Symposium on Foundations of Computer Science (FOCS 2016), 2016, pp. 438–446.
- [32] P. McCULLAGH, *Tensor Methods in Statistics*, Courier Dover, Mineola, NY, 2018.
- [33] T. ROYEN, *A simple proof of the Gaussian correlation conjecture extended to multivariate gamma distributions*, Far East J. Theor. Stat., 48 (2014), pp. 139–145, <http://www.pphmj.com/abstract/8713.htm>.
- [34] Z. ŠIDÁK, *Rectangular confidence regions for the means of multivariate normal distributions*, J. Amer. Statist. Assoc., 62 (1967), pp. 626–633, <https://doi.org/10.1080/01621459.1967.10482935>.
- [35] G. W. STEWART, *Perturbation Theory for the Singular Value Decomposition*, technical report, University of Maryland Institute for Advanced Computer Studies, College Park, MD, 1990.
- [36] G. W. STEWART, *Perturbation theory for the singular value decomposition*, in SVD and Signal Processing, II: Algorithms, Analysis and Applications, R. J. Vacarro, ed., Elsevier, Amsterdam, 1991, pp. 99–109.