

UC Irvine

UC Irvine Previously Published Works

Title

Analyzing Tensor Power Method Dynamics in Overcomplete Regime

Permalink

<https://escholarship.org/uc/item/58r0q4hw>

Authors

Anandkumar, Animashree

Ge, Rong

Janzamin, Majid

Publication Date

2017

Peer reviewed

Analyzing Tensor Power Method Dynamics in Overcomplete Regime

Anima Anandkumar* Rong Ge† Majid Janzamin‡

September 16, 2015

Abstract

We present a novel analysis of the dynamics of tensor power iterations in the overcomplete regime where the tensor CP rank is larger than the input dimension. Finding the CP decomposition of an overcomplete tensor is NP-hard in general. We consider the case where the tensor components are randomly drawn, and show that the simple power iteration recovers the components with bounded error under mild initialization conditions. We apply our analysis to unsupervised learning of latent variable models, such as multi-view mixture models and spherical Gaussian mixtures. Given the third order moment tensor, we learn the parameters using tensor power iterations. We prove it can correctly learn the model parameters when the number of hidden components k is much larger than the data dimension d , up to $k = o(d^{1.5})$. We initialize the power iterations with data samples and prove its success under mild conditions on the signal-to-noise ratio of the samples. Our analysis significantly expands the class of latent variable models where spectral methods are applicable. Our analysis also deals with noise in the input tensor leading to sample complexity result in the application to learning latent variable models.

Keywords: tensor decomposition, tensor power iteration, overcomplete representation, unsupervised learning, latent variable models.

1 Introduction

CANDECOMP/PARAFAC (CP) decomposition of a symmetric tensor $T \in \mathbb{R}^{d \times d \times d}$ is the process of decomposing it into a succinct sum of rank-one tensors, given by

$$T = \sum_{j \in [k]} \lambda_j a_j \otimes a_j \otimes a_j, \quad \lambda_j \in \mathbb{R}, \quad a_j \in \mathbb{R}^d, \quad (1)$$

where \otimes denotes the outer product. The minimum k for which the tensor can be decomposed in the above form is called the (symmetric) tensor rank. *Tensor power iteration* is a simple, popular and

*University of California, Irvine. Email: a.anandkumar@uci.edu

†Duke University. Email: rongge@cs.duke.edu

‡University of California, Irvine. Email: mjanzami@uci.edu

efficient method for recovering the tensor rank-one components a_j 's. The tensor power iteration is given by

$$x \leftarrow \frac{T(I, x, x)}{\|T(I, x, x)\|}, \quad (2)$$

where

$$T(I, x, x) := \sum_{j,l \in [d]} x_j x_l T(:, j, l) \in \mathbb{R}^d$$

is a *multilinear* combination of tensor *fibers*, and $\|\cdot\|$ is the ℓ_2 norm operator. See Section 1.3 for an overview of tensor notations and preliminaries.

The tensor power iteration is a generalization of matrix power iteration: for matrix $M \in \mathbb{R}^{d \times d}$, the power iteration is given by $x \leftarrow Mx/\|Mx\|$. Dynamics and convergence properties of matrix power iterations are well understood (Horn and Johnson, 2012). On the other hand, a theoretical understanding of tensor power iterations is much more limited. Tensor power iteration can be viewed as a *gradient descent* step (with infinite step size), corresponding to the problem of finding the best rank-1 approximation of the input tensor T (Anandkumar et al., 2014c). This optimization problem is non-convex. Unlike the matrix case, where the number of isolated stationary points of power iteration is at most the dimension (given by eigenvectors corresponding to unique eigenvalues), in the tensor case, the number of stationary points is, in fact, exponential in the input dimension (Cartwright and Sturmfels, 2013). This makes the analysis of tensor power iteration far more challenging.

Despite the above challenges, many advances have been made in understanding the tensor power iterations in specific regimes. When the components a_j 's are orthogonal to one another, it is known that there are no spurious local optima for tensor power iterations, and the only stable fixed points correspond to the true a_j 's (Zhang and Golub, 2001; Anandkumar et al., 2014c). Any tensor with linearly independent components a_j 's can be orthogonalized, via an invertible transformation (whitening) and thus, its components can be recovered efficiently. A careful perturbation analysis in this setting was carried out in Anandkumar et al. (2014c).

The framework in Anandkumar et al. (2014c) is however not applicable in the overcomplete setting, where the tensor rank k exceeds the dimension d . Such overcomplete tensors cannot be orthogonalized and finding guaranteed decomposition is a challenging open problem. It is known that finding CP tensor decomposition is NP-hard (Hillar and Lim, 2013). In this paper, we make significant headway in showing that the simple power iterations can recover the components in the overcomplete regime under a set of mild conditions on the components a_j 's.

Overcomplete tensors also arise in many machine learning applications such as moments of many latent variable models, e.g., multiview mixtures, independent component Analysis (ICA), and sparse coding models, where the number of hidden variables exceeds the input dimensions (Anandkumar et al., 2015). Overcomplete models often have impressive empirical performance (Coates et al., 2011), and can provide greater flexibility in modeling, and are more robust to noise (Lewicki and Sejnowski, 2000). By studying algorithms for overcomplete tensor decomposition, we expand the class of models that can be learnt efficiently using simple spectral methods such as tensor power iterations. Note there are other algorithms for decomposing overcomplete tensors (De Lathauwer et al., 2007; Goyal et al., 2013; Bhaskara et al., 2013), but they all require tensors of at least 4-th order and require large computational complexity. Ge and Ma (2015) works for 3rd order tensor but requires quasi-polynomial time. The main contribution of this paper is an analysis for the practical power method in the overcomplete regime.

1.1 Summary of results

We analyze the dynamics of third order tensor power iterations in the overcomplete regime. We assume that the tensor components a_j 's are randomly drawn from the unit sphere. Since general tensor decomposition is challenging in the overcomplete regime, we argue that this is a natural first step to consider for tractable recovery.

We characterize the basin of attraction for the local optima near the rank-one components a_j 's. We show that under mild initialization condition, there is fast convergence to these local optima in $O(\log \log d)$ iterations (i.e., quadratic convergence as opposed to linear convergence in case of matrices). This result is the core technical analysis of this paper stated in the following theorem.

Theorem 1 (Dynamics of tensor power iteration). *Consider tensor $\hat{T} = T + E$ such that exact tensor T has rank- k decomposition in (1) with rank-one components $a_j \in \mathbb{R}^d, j \in [k]$ being uniformly i.i.d. drawn from the unit d -dimensional sphere, and the ratio of maximum and minimum (in absolute value) weights λ_j 's being constant. In addition, suppose the perturbation tensor E has bounded norm as*

$$\|E\| \leq \epsilon \frac{\sqrt{k}}{d}, \quad \text{where } \epsilon < o\left(\frac{\sqrt{k}}{d}\right). \quad (3)$$

Let tensor rank $k = o(d^{1.5})$, and the unit-norm initial vector $x^{(1)}$ satisfy the correlation bound

$$|\langle x^{(1)}, a_j \rangle| \geq d^\beta \frac{\sqrt{k}}{d}, \quad (4)$$

w.r.t. some true component $a_j, j \in [k]$, for some constant $\beta > 0$. After $N = \Theta(\log \log d)$ iterations, the tensor power iteration in (2) outputs a vector having w.h.p. a constant correlation with the true component a_j as $|\langle x^{(N+1)}, a_j \rangle| \geq 1 - \gamma$, for any fixed constant $\gamma > 0$.

As a corollary, this result can be used for learning latent variable models such as multiview mixtures. We show that the above initialization condition is satisfied using a sample with mild signal-to-noise ratio; see Section 2 for more details on this.

The above result is a significant improvement over the recent analysis by Anandkumar et al. (2015, 2014a,b) for overcomplete tensor decomposition. In these works, it is required for the initialization vectors to have a constant amount of correlation with the true a_j 's. However, obtaining such strong initializations is usually not realistic in practice. On the other hand, the initialization condition in (4) is mild, and decaying even when the rank k is significantly larger than dimension d ; up to $k = o(d^{1.5})$. In learning the mixture model, such initialization vectors can be obtained as samples from the mixture model, even when there is a large amount of noise. Given this improvement, we combine our analysis in Theorem 1, and the guarantees in (Anandkumar et al., 2014a), proving that the model parameters can be recovered consistently.

A detailed proof outline for Theorem 1 is provided in Section 3.1. Under the random assumption, it is not hard to show that the first iteration of tensor power update makes progress. However, after the first iteration, the input vector and the tensor components are no longer *independent* of each other. Therefore, we cannot directly repeat the same argument for the second step.

How do we analyze the second step even though the vector and tensor components are correlated? The main intuition is to characterize the dependency between the vector and the tensor components, and show that there is still enough randomness left for us to repeat the argument. This idea was inspired by the analysis of Approximate Message Passing (AMP) algorithms (Bayati and Montanari, 2010). However, our analysis here is very different in several key

aspects: 1) In approximate message passing, typically the analysis works in the *large system limit*, where the number of iterations is fixed and the dimension goes to infinity. Here we can handle a superconstant number of iterations $O(\log \log d)$, even for finite d ; 2) Usually k is assumed to be a constant factor times d in the AMP-like analysis, while here we allow them to be polynomially related.

1.2 Related work

Tensor decomposition for learning latent variable models: In the introduction, some related works are mentioned which study the theoretical and practical aspects of spectral techniques for learning latent variable models. Among them, Anandkumar et al. (2014c) provide the analysis of tensor power iteration for learning several latent variable models in the undercomplete regime. Anandkumar et al. (2014a) provide the analysis in the overcomplete regime and Anandkumar et al. (2014b) provide tensor concentration bounds and apply the analysis in (Anandkumar et al., 2014a) to learning LVMs proposing tight sample complexity guarantees.

Learning mixture of Gaussians: Here, we provide a subset of related works studying learning mixture of Gaussians which are more comparable with our result. For a more detailed list of these works, see Anandkumar et al. (2014c); Hsu and Kakade (2013). The problem of learning mixture of Gaussians dates back to the work by Pearson (1895). They propose a moment-based technique that involves solving systems of multivariate polynomials which is in general challenging in both computational and statistical sense. Recently, lots of studies on learning Gaussian mixture models have been done improving both aspects which can be divided to two main classes: distance-based and spectral methods.

Distance-based methods impose separation condition on the mean vectors showing that under enough separation the parameters can be estimated. Among such approaches, we can mention Dasgupta (1999); Vempala and Wang (2002); Arora and Kannan (2005). As discussed in the summary of results, these results work even if $k > d^{1.5}$ as long as the separation condition between means is satisfied, but our work can tolerate higher level of noise in the regime of $k = o(d^{1.5})$ with polynomial computational complexity. The guarantees in (Vempala and Wang, 2002) also work in the high noise regime but need higher computational complexity as polynomial in $k^{O(k)}$ and d .

In the spectral approaches, the observed moments are constructed and the spectral decomposition of the observed moments are performed to recover the parameters (Kalai et al., 2010; Anandkumar et al., 2012, 2014b). Kalai et al. (2010) analyze the problem of learning mixture of two general Gaussians and provide algorithm with high order polynomial sample and computational complexity. Note that in general, the complexity of such methods grow exponentially with the number of components without further assumptions (Moitra and Valiant, 2010). Hsu and Kakade (2013) provide a spectral algorithm under non-degeneracy conditions on the mean vectors and providing guarantees with polynomial sample complexity depending on the condition number of the moment matrices. Anandkumar et al. (2014b) perform tensor power iteration on the third order moment tensor to recover the mean vectors in the overcomplete regime as long as $k = o(d^{1.5})$, but need very good initialization vector having constant correlation with the true mean vector. Here, we improve the correlation level required for convergence.

1.3 Notation and tensor preliminaries

Let $[k] := \{1, 2, \dots, k\}$, and $\|v\|$ denote the ℓ_2 norm of vector v . We use \tilde{O} and $\tilde{\Omega}$ to hide polylog factors in asymptotic notations O and Ω , respectively.

Tensor preliminaries: A real p -th order tensor $T \in \otimes^p \mathbb{R}^d$ is a member of the outer product of Euclidean spaces \mathbb{R}^d . The different dimensions of the tensor are referred to as *modes*. For instance, for a matrix, the first mode refers to columns and the second mode refers to rows. In addition, *fibers* are higher order analogues of matrix rows and columns. A fiber is obtained by fixing all but one of the indices of the tensor (and is arranged as a column vector). For example, for a third order tensor $T \in \mathbb{R}^{d \times d \times d}$, the mode-1 fiber is given by $T(:, j, l)$. Similarly, *slices* are obtained by fixing all but two of the indices of the tensor. For example, for the third order tensor T , the slices along 3rd mode are given by $T(:, :, l)$.

We view a tensor $T \in \mathbb{R}^{d \times d \times d}$ as a multilinear form. In particular, for vectors $u, v, w \in \mathbb{R}^d$, we have¹

$$T(I, v, w) := \sum_{j, l \in [d]} v_j w_l T(:, j, l) \in \mathbb{R}^d, \quad (5)$$

which is a multilinear combination of the tensor mode-1 fibers. Similarly $T(u, v, w) \in \mathbb{R}$ is a multilinear combination of the tensor entries, and $T(I, I, w) \in \mathbb{R}^{d \times d}$ is a linear combination of the tensor slices.

A 3rd order tensor $T \in \mathbb{R}^{d \times d \times d}$ is said to be rank-1 if it can be written in the form

$$T = \lambda \cdot a \otimes b \otimes c \Leftrightarrow T(i, j, l) = \lambda \cdot a(i) \cdot b(j) \cdot c(l), \quad (6)$$

where notation \otimes represents the *outer product* and $a, b, c \in \mathbb{R}^d$ are unit vectors. A tensor $T \in \mathbb{R}^{d \times d \times d}$ is said to have a CP *rank* at most k if it can be written as the sum of k rank-1 tensors as

$$T = \sum_{i \in [k]} \lambda_i a_i \otimes b_i \otimes c_i, \quad \lambda_i \in \mathbb{R}, \quad a_i, b_i, c_i \in \mathbb{R}^d. \quad (7)$$

In the rest of the paper, Section 2 describes how to apply our tensor results to learning multiview mixture models. Section 3 illustrates the proof ideas, with more details in the Appendix. Finally we conclude in Section 4.

2 Learning multiview mixture model through tensor methods

We proposed our main technical result in Section 1.1 providing convergence guarantees for the tensor power iterations given mild initialization conditions in the overcomplete regime; see Theorem 1. Along this result we provide the application to learning multiview mixtures model in Theorem 2. In this section, we briefly introduce the tensor decomposition framework as the learning algorithm and then state the learning guarantees with more details and remarks.

¹Compare with the matrix case where for $M \in \mathbb{R}^{d \times d}$, we have $M(I, u) = Mu := \sum_{j \in [d]} u_j M(:, j)$.

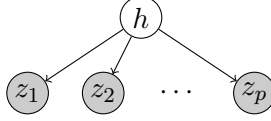


Figure 1: Multiview mixture model.

2.1 Multiview mixture model

Consider an exchangeable multiview mixture model with k components and $p \geq 3$ views; see Figure 1. Suppose that hidden variable h is a discrete categorical random variable taking one of the k states. It is convenient to represent it by basis vectors such that

$$h = e_j \in \mathbb{R}^k \quad \text{if and only if} \quad \text{it takes the } j\text{-th state.}$$

Note that $e_j \in \mathbb{R}^k$ denotes the j -th basis vector in the k -dimensional space. The prior probability for each hidden state is also $\Pr[h = e_j] = \lambda_j, j \in [k]$. For simplicity, in this paper we assume all the λ_i 's are the same. However, similar argument works even when the ratio of maximum and minimum prior probabilities $\lambda_{\max}/\lambda_{\min}$ is bounded by some constant.

The variables (views) $z_l \in \mathbb{R}^d$ are related to the hidden state through *factor matrix* $A \in \mathbb{R}^{d \times k}$ such that

$$z_l = Ah + \eta_l, \quad l \in [p],$$

where zero-mean noise vectors $\eta_l \in \mathbb{R}^d$ are independent of each other and the hidden state h . Given this, the variables (views) $z_l \in \mathbb{R}^d$ are conditionally independent given the latent variable h , and the conditional means are $\mathbb{E}[z_l|h = e_j] = a_j$, where $a_j \in \mathbb{R}^d$ denotes the j -th column of factor matrix $A = [a_1 \cdots a_k] \in \mathbb{R}^{d \times k}$. In addition, the above properties imply that the order of observations z_l do not matter and the model is *exchangeable*. The goal of the learning problem is to recover the parameters of the model (factor matrix) A given observations.

For this model, the third order² observed moment has the form (Anandkumar et al., 2014c)

$$\mathbb{E}[z_1 \otimes z_2 \otimes z_3] = \sum_{j \in [k]} \lambda_j a_j \otimes a_j \otimes a_j. \quad (8)$$

Hence, given third order observed moment, the unsupervised learning problem (recovering factor matrix A) reduces to computing a tensor decomposition as in (8).

2.2 Tensor decomposition algorithm

The algorithm for unsupervised learning of multiview mixture model is based on tensor decomposition techniques provided in Algorithm 1. The main step in (9) performs *tensor power iteration*³; see (5) for the multilinear form definition. After running the algorithm for all different initialization vectors, the clustering process from Anandkumar et al. (2014a) ensures that the best converged vectors are returned as the estimation of true components a_j .

²It is enough to form the third order moment for our learning purpose.

³This is the generalization of matrix power iteration to 3rd order tensors.

Algorithm 1 Learning multiview mixture model via tensor power iterations

Require: 1) Third order moment tensor $T \in \mathbb{R}^{d \times d \times d}$ in (8), 2) n samples of z_1 in multiview mixture model as $z_1^{(\tau)}, \tau \in [n]$, and 3) number of iterations N .

- 1: **for** $\tau = 1$ **to** n **do**
- 2: **Initialize** unit vectors $x_\tau^{(1)} \leftarrow z_1^{(\tau)} / \|z_1^{(\tau)}\|$.
- 3: **for** $t = 1$ **to** N **do**
- 4: Tensor power updates (see (5) for the definition of the multilinear form):

$$x_\tau^{(t+1)} = \frac{T(I, x_\tau^{(t)}, x_\tau^{(t)})}{\|T(I, x_\tau^{(t)}, x_\tau^{(t)})\|}, \quad (9)$$

- 5: **end for**
 - 6: **end for**
 - 7: **return** the output of Procedure 2 with input $\{x_\tau^{(N+1)} : \tau \in [n]\}$ as estimates x_j .
-

Procedure 2 Clustering process (Anandkumar et al., 2014a)

Require: Tensor $T \in \mathbb{R}^{d \times d \times d}$, set $S := \{x_\tau^{(N+1)} : \tau \in [n]\}$, parameter ν .

- 1: **while** S is not empty **do**
 - 2: Choose $x \in S$ which maximizes $|T(x, x, x)|$.
 - 3: Do N more iterations of (9) starting from x .
 - 4: **Output** the result of iterations denoted by \hat{x} .
 - 5: Remove all the $x \in S$ with $|\langle x, \hat{x} \rangle| > \nu/2$.
 - 6: **end while**
-

2.3 Learning guarantees

We assume a Gaussian prior on the mean vectors, i.e., the vectors $a_j \sim \mathcal{N}(0, I_d/d)$, $j \in [k]$ are i.i.d. drawn from a standard multivariate Gaussian distribution with unit expected square norm. Note that in the high dimension (growing d), this assumption is the same as uniformly drawing from unit sphere since the norm of vector concentrates in the high dimension and there is no need for normalization. Even though we impose a prior distribution, we do not use a MAP estimator, since the corresponding optimization is NP-hard. Instead, we learn the model parameters through decomposition of the third order moments through tensor power iterations. The assumption of a Gaussian prior is standard in machine learning applications. We impose it here for tractable analysis of power iteration dynamics. Such Gaussian assumptions have been used before for analysis of other iterative methods such as approximate message passing algorithms, and there are evidences that similar results hold for more general distributions; see (Bayati and Montanari, 2010) and references there.

As explained in the previous sections, we use tensor power method to learn the components a_j 's, and the method is initialized with observed samples z_i . Intuitively, this initialization is useful since $z_i = Ah + \eta_i$ is a perturbed version of desired parameter a_j (when $h = e_j$). Thus, we present the result in terms of the signal-to-noise (SNR) ratio which is the expected norm of signal a_j (which is one here) divided by the expected norm of noise η_i , i.e., the SNR in the i -th sample

$z_i = a_j + \eta_i$ (assumed $h = e_j$) is defined as $\text{SNR} := \mathbb{E}[\|a_j\|]/\mathbb{E}[\|\eta_i\|]$. This specifies how much noise the initialization vector z_i can tolerate in order to ensure the convergence of tensor power iteration to a desired local optimum. We now propose the conditions required for recovery guarantees, and state a brief explanation of them.

Conditions for Theorems 2 and 3:

- Rank condition: $k \leq o(d^{1.5})$.
- The columns of A are uniformly i.i.d. drawn from unit d -dimensional sphere.
- The noise vectors $\eta_l, l \in [3]$, are independent of matrix A and each other. In addition, the signal-to-noise ratio (SNR) is w.h.p. bounded as

$$\text{SNR} \geq \Omega \left(\frac{\sqrt{\max\{k, d\}}}{d^{1-\beta}} \right),$$

for some $\beta \geq (\log d)^{-c}$ for universal constant $c > 0$.

The rank condition bounds the level of overcompleteness for which the recovery guarantees are satisfied. The random assumption on the columns of A are crucial for analyzing the dynamics of tensor power iteration. We use it to argue there exists enough randomness left in the components after conditioning on the previous iterations; see Section 3.1 for the details. The bound on the SNR is required to make sure the given sample used for initialization is close enough to the corresponding mean vector. This ensures that the initial vector is inside the basin-of-attraction of the corresponding component, and hence, the convergence to the mean vector can be guaranteed. Under these assumptions we have.

Theorem 2 (Learning multiview mixture model given exact tensor: closeness to single columns). *Consider a multiview mixture model (or a spherical Gaussian mixture) in the above setting with k components in d dimensions. If the above conditions hold, then the tensor power iteration converges to a vector close to one of the true mean vectors a_j 's (having constant correlation).*

In particular, for mildly overcomplete models, where $k = \alpha d$ for some constant $\alpha > 1$, the signal-to-noise ratio (SNR) is as low as $\Omega(d^{-1/2+\epsilon})$, for any $\epsilon > 0$. Thus, we can learn mixture models with a high level of noise. In general, we establish how the required noise level scales with the number of hidden components k , as long as $k = o(d^{1.5})$.

The above theorem states convergence to desired local optima which are close to true components a_j 's. In Theorem 3, we show that we can sharpen the above result, by jointly iterating over the recovered vectors, and consistently recover the components a_j 's. This result also uses the analysis from Anandkumar et al. (2015).

Theorem 3 (Learning multiview mixture model given exact tensor: recovering the whole factor matrix). *Assume the above conditions hold. The initialization of power iteration is performed by samples of z_1 in multiview mixture model. Suppose the tensor power iterations is at least initialized once for each $a_j, j \in [k]$ such that $z_1 = a_j + \eta_1$.⁴ Then by using the exact 3rd order moment tensor*

⁴Note that this happens for component j with high probability when the number of initializations is proportional to inverse prior probability corresponding to that mixture.

in (8) as input, the tensor decomposition algorithm outputs an estimate \hat{A} satisfying w.h.p. (over the randomness of the components a_j 's)

$$\left\| \hat{A} - A \right\|_F \leq \epsilon,$$

where the number of iterations of the algorithm is $N = \Theta \left(\log \left(\frac{1}{\epsilon} \right) + \log \log d \right)$.

The above theorems assume the exact third order tensor is given to the algorithm. We provide the results given empirical tensor in Section 2.3.1.

Learning spherical Gaussian mixtures: Consider a mixture of k different Gaussian vectors with spherical covariance. Let $a_j \in \mathbb{R}^d, j \in [k]$ denote the mean vectors and the covariance matrices are $\sigma^2 I$. Assuming the parameter σ is known, the modified third order observed moment

$$M_3 := \mathbb{E}[z \otimes z \otimes z] - \sigma^2 \sum_{i \in [d]} (\mathbb{E}[z] \otimes e_i \otimes e_i + e_i \otimes \mathbb{E}[z] \otimes e_i + e_i \otimes e_i \otimes \mathbb{E}[z])$$

has the tensor decomposition form (Hsu and Kakade, 2012)

$$M_3 = \sum_{j \in [k]} \lambda_j a_j \otimes a_j \otimes a_j,$$

where λ_j is the probability of drawing j -th Gaussian mixture. The above guarantees can be applied to learning mean vectors a_j in this model with the additional property that the noise is spherical Gaussian.

Learning multiview mixture model with distinct factor matrices: Consider the multiview mixture model with different factor matrices where the first three views are related to the hidden state as

$$z_1 = Ah + \eta_1, \quad z_2 = Bh + \eta_2, \quad z_3 = Ch + \eta_3.$$

Then, the guarantees in the above theorem can be extended to recovering the columns of all three factor matrices A , B , and C with appropriate modifications in the power iteration algorithm as follows. First the update formula (9) is changed as

$$x_{1,\tau}^{(t+1)} = \frac{T \left(I, x_{2,\tau}^{(t)}, x_{3,\tau}^{(t)} \right)}{\left\| T \left(I, x_{2,\tau}^{(t)}, x_{3,\tau}^{(t)} \right) \right\|}, \quad x_{2,\tau}^{(t+1)} = \frac{T \left(x_{1,\tau}^{(t)}, I, x_{3,\tau}^{(t)} \right)}{\left\| T \left(x_{1,\tau}^{(t)}, I, x_{3,\tau}^{(t)} \right) \right\|}, \quad x_{3,\tau}^{(t+1)} = \frac{T \left(x_{1,\tau}^{(t)}, x_{2,\tau}^{(t)}, I \right)}{\left\| T \left(x_{1,\tau}^{(t)}, x_{2,\tau}^{(t)}, I \right) \right\|},$$

which is the alternating asymmetric version of symmetric power iteration in (9). Here, we alternate among different modes of the tensor. In addition, the initialization for each mode of the tensor is appropriately performed with the samples corresponding to that mode. Note that the analysis still works in the asymmetric version since there exists even more independence relationships through the iterations of the power update because of introducing new random matrices B and C .

2.3.1 Sample complexity analysis

In the previous section, we assumed the exact third order tensor in (8) is given to the tensor decomposition Algorithm 1. We now estimate the tensor given n samples $z_1^{(i)}, z_2^{(i)}, z_3^{(i)}, i \in [n]$, as

$$\hat{T} = \frac{1}{n} \sum_{i \in [n]} z_1^{(i)} \otimes z_2^{(i)} \otimes z_3^{(i)}. \quad (10)$$

For the multiview mixture model introduced in Section 2.1, let the noise vector η_l be spherical, and ζ^2 denote the variance of each entry of noise vector. We now provide the following recovery guarantees.

Additional conditions for Theorem 4:

- Let $E_1 := [\eta_1^{(1)}, \eta_1^{(2)}, \dots, \eta_1^{(n)}] \in \mathbb{R}^{d \times n}$, where $\eta_1^{(i)} \in \mathbb{R}^d$ is the i -th sample of noise vector η_1 . These noise matrices satisfy the following *RIP property* which is adapted from Candes and Tao (2006). Matrix $E_1 \in \mathbb{R}^{d \times n}$ satisfies a weak RIP condition such that for any subset of $O\left(\frac{d}{\log^2 d}\right)$ number of columns, the spectral norm of E_1 restricted to those columns is bounded by 2. The same condition is satisfied for similarly defined noise matrices E_2 and E_3 .
- The number of samples n satisfies lower bound such that

$$\zeta \left(\frac{\sqrt{d}}{n} + \sqrt{\lambda_{\max} \frac{d}{n}} \right) + \zeta^2 \left(\frac{d}{n} + \sqrt{\lambda_{\max} \frac{d^{1.5}}{n}} \right) + \zeta^3 \left(\frac{d^{1.5}}{n} + \sqrt{\frac{d}{n}} \right) \leq \min \left\{ \epsilon \frac{\sqrt{k}}{d}, \tilde{O}(\lambda_{\min}) \right\}, \quad (11)$$

where $\epsilon < o(\sqrt{k}/d)$.

Theorem 4 (Learning multiview mixture model given empirical tensor). *Consider the empirical tensor in (10) as the input to tensor decomposition Algorithm 1. Suppose the above additional conditions are also satisfied. Then, the same guarantees as in Theorem 2 hold. In addition, the same guarantees as in Theorem 3 also hold with the recovery bound changed as*

$$\|\hat{A} - A\|_F \leq \tilde{O} \left(\frac{\sqrt{k} \cdot \|E\|}{\lambda_{\min}} \right),$$

where E denotes the perturbation tensor originated from empirical estimation in (10), and its spectral norm $\|E\|$ is bounded by the LHS of (11).

Proof: The above sample complexity result is proved by using the tensor concentration bound in Theorem 1 of Anandkumar et al. (2014b) applied to our noisy analysis of tensor power dynamics in Theorem 1; see Equation (3). The additional bound on sample complexity and final recovery error on $\|\hat{A} - A\|_F$ is also from Theorem 1 of Anandkumar et al. (2015). \square

3 Proof Outline

Our main technical result is the analysis of third order tensor power iteration provided in Theorem 1 which also allows to tolerate some amount of noise in the input tensor. We analyze the noiseless and

noisy settings in different ways. We basically first prove the result for the noiseless setting where the input tensor has an exact rank- k decomposition in (1). When the noise is also considered, we show that the contribution of noise in the analysis is dominated by the main signal, and thus, the same result still holds. For the rest of this section we focus on the noiseless setting, while we discuss the proof ideas for the noisy case in Section 3.2.

We first discuss the proof of Theorem 3 which involves two phases. In the first phase, we show that under certain small amount of correlation (see (13)) between the initial vector and the true component, the power iteration in (2) converges to some vector which has constant correlation with the true component. This result is the core technical analysis of this paper which is provided in Lemma 5. In the second phase, we incorporate the result of Anandkumar et al. (2014a) which guarantees the approximate convergence of power iteration given initial vector having constant correlation with the true component. This is stated in Lemma 6.

To simplify the notation, we consider the tensor⁵

$$T = \sum_{j \in [k]} a_j \otimes a_j \otimes a_j, \quad a_j \sim \mathcal{N}(0, \frac{1}{d} I_d). \quad (12)$$

Notice that this is exactly proportional to the 3rd order moment tensor of the multiview mixture model in (8).

The following lemma is restatement of Theorem 1 in the noiseless setting.

Lemma 5 (Dynamics of tensor power iteration, phase 1). *Consider the rank- k tensor T of the form in (12). Let tensor rank $k = o(d^{1.5})$, and the unit-norm initial vector $x^{(1)}$ satisfies the correlation bound*

$$|\langle x^{(1)}, a_j \rangle| \geq d^\beta \frac{\sqrt{k}}{d}, \quad (13)$$

w.r.t. some true component $a_j, j \in [k]$, for some $\beta > (\log d)^{-c}$ for some universal constant $c > 0$. After $N = \Theta(\log \log d)$ iterations, the tensor power iteration in (2) outputs a vector having w.h.p. a constant correlation with the true component a_j as

$$|\langle x^{(N+1)}, a_j \rangle| \geq 1 - \gamma,$$

for any fixed constant $\gamma > 0$.

The proof outline of above lemma is provided in Section 3.1.

Lemma 6 (Dynamics of tensor power iteration, phase 2 (Anandkumar et al., 2014a)). *Consider the rank- k tensor T of the form in (12) with rank condition $k \leq o(d^{1.5})$. Let the initial vectors $x_j^{(1)}$ satisfy the constant correlation bound*

$$|\langle x_j^{(1)}, a_j \rangle| \geq 1 - \gamma_j,$$

w.r.t. true components $a_j, j \in [k]$, for some constants $\gamma_j > 0$. Let the output of the tensor power update⁶ in (2) applied to all these different initialization vectors after $N = \Theta(\log \frac{1}{\epsilon})$ iterations be

⁵In the analysis, we assume that all the weights are equal to one which can be generalized to the case when the ratio of maximum and minimum weights (in absolute value) are constant.

⁶This result also needs an additional step of coordinate descent iterations since the true components are not the fixed points of power iteration; see Anandkumar et al. (2014a) for the details.

stacked in matrix \hat{A} . Then, we have w.h.p.⁷

$$\left\| \hat{A} - A \right\|_F \leq \epsilon.$$

Given the above two lemmas, the learning result in Theorem 3 is directly proved.

Proof of Theorem 3: The result is proved by combining Lemma 5 and Lemma 6. Note that the initialization condition in (4) is w.h.p. satisfied given the SNR bound assumed. \square

3.1 Proof outline of Lemma 5 (noiseless case of Theorem 1)

First step: We first intuitively show the first step of the algorithm makes progress. Suppose the tensor is $T = \sum_{j \in [k]} a_j \otimes a_j \otimes a_j$, and the initial vector x has correlation $|\langle x, a_1 \rangle| \geq d^\beta \frac{\sqrt{k}}{d}$ with the first component. The result of the first iteration is the normalized version of the following vector:

$$\tilde{x} = \sum_{j \in [k]} \langle a_j, x \rangle^2 a_j.$$

Intuitively, this vector should have roughly $\langle a_1, \tilde{x} \rangle = d^{2\beta} \frac{k}{d^2}$ correlation with a_1 (as the other terms are random they don't contribute much). On the other hand, the norm of this vector is roughly $O(\sqrt{k}/d)$: this is because $\langle a_j, x \rangle^2$ for $j \neq 1$ is roughly⁸ $1/d$, and the sum of k random vectors with length $1/d$ will have length roughly $O(\sqrt{k}/d)$. These arguments can be made precise showing the normalized version $\tilde{x}/\|\tilde{x}\|$ has correlation $d^{2\beta} \frac{\sqrt{k}}{d}$ with a_1 ensuring progress in the first step.

Going forward: As we explained, the basic idea behind proving Lemma 5 is to characterize the conditional distribution of random Gaussian tensor components a_j 's given previous iterations. In particular, we show that the residual independent randomness left in these conditional distributions is large enough and we can exploit it to obtain tighter concentration bounds throughout the analysis of the iterations. The Gaussian assumption on the components, and small enough number of iterations are crucial in this argument.

Notations: For two vectors $u, v \in \mathbb{R}^k$, the Hadamard product denoted by $*$ is defined as the entry-wise multiplication of vectors, i.e., $(u * v)_j := u_j v_j$ for $j \in [k]$. For a matrix A , let $P_{\perp A}$ denote the projection operator to the subspace orthogonal to column span of A . For a subspace R , let R^\perp denote the space orthogonal to it. Therefore, for a subspace R , the projection operator on the subspace orthogonal to R is equivalently denoted by P_{R^\perp} or $P_{\perp R}$. For a random matrix D , let $D|\{u = Dv\}$ denote the conditional distribution of D given linear constraints $u = Dv$.

Lemma 5 involves analyzing the dynamics of power iteration in (2) for 3rd order rank- k tensors. For the rank- k tensor in (12), the power iterative form $x \leftarrow \frac{T(I, x, x)}{\|T(I, x, x)\|}$ can be written as

$$x^{(t+1)} = \frac{A (A^\top x^{(t)})^{*2}}{\left\| A (A^\top x^{(t)})^{*2} \right\|}, \quad (14)$$

⁷Anandkumar et al. (2014a) recover the vector up to sign since they work in the asymmetric case. In symmetric case it is easy to resolve sign ambiguity issue.

⁸The correlation between two unit Gaussian vectors in d dimensions is roughly $1/\sqrt{d}$.

where the multilinear form in (5) is used. Here, $A = [a_1 \cdots a_k] \in \mathbb{R}^{d \times k}$ denotes the factor matrix, and for vector $y \in \mathbb{R}^k$, $y^{*2} := y * y \in \mathbb{R}^k$ represents the element-wise square of entries of y .

We consider the case where $a_i \sim \mathcal{N}(0, \frac{1}{d}I)$ are i.i.d. drawn and we analyze the evolution of the dynamics of the power update. As explained earlier, for a given initialization $x^{(1)}$, the update in the first step can be analyzed easily since A is independent of $x^{(1)}$. However, in subsequent steps, the updates $x^{(t)}$ are dependent on A , and it is no longer clear how to provide a tight bound on the evolution of $x^{(t)}$. In this work, we provide a careful analysis by controlling the amount of ‘‘correlation build-up’’ by exploiting the structure of Gaussian matrices under linear constraints. This enables us to provide better guarantees for matrix A with Gaussian entries compared to general matrices A .

Intermediate update steps and variables: Before we proceed, we need to break down power update in (2) and introduce some intermediate update steps and variables as follows. Recall that $x^{(1)} \in \mathbb{R}^d$ denotes the initialization vector. Without loss of generality, let us analyze the convergence of power update to first component of rank- k tensor T denoted by a_1 . Hence, let the first entry of $x^{(1)}$ denoted by $x_1^{(1)}$ be the maximum entry (in absolute value) of $x^{(1)}$, i.e., $x_1^{(1)} = \|x^{(1)}\|_\infty$. Let $B := [a_2 \ a_3 \ \cdots \ a_k] \in \mathbb{R}^{d \times (k-1)}$, and therefore $A = [a_1|B]$. We break the power update formula in (2) into a few steps by introducing intermediate variables $y^{(t)} \in \mathbb{R}^k$ and $\tilde{x}^{(t+1)} \in \mathbb{R}^d$ as

$$y^{(t)} := A^\top x^{(t)}, \quad \tilde{x}^{(t+1)} := A(y^{(t)})^{*2}.$$

Note that $\tilde{x}^{(t+1)}$ is the unnormalized version of $x^{(t+1)} := \tilde{x}^{(t+1)} / \|\tilde{x}^{(t+1)}\|$, i.e., $\tilde{x}^{(t+1)} := T(I, x^{(t)}, x^{(t)})$. Thus, we need to jointly analyze the dynamics of all variables $x^{(t)}$, $y^{(t)}$ and $(y^{(t)})^{*2}$. Define

$$X^{[t]} := [x^{(1)} | \cdots | x^{(t)}], \quad Y^{[t]} := [y^{(1)} | \cdots | y^{(t)}].$$

Matrix B is randomly drawn with i.i.d. Gaussian entries $B_{ij} \sim \mathcal{N}(0, \frac{1}{d})$. As the iterations proceed, we consider the following conditional distributions

$$B^{(t,1)} := B|\{X^{[t]}, Y^{[t]}\}, \quad B^{(t,2)} := B|\{X^{[t+1]}, Y^{[t]}\}. \quad (15)$$

Thus, $B^{(t,1)}$ is the conditional distribution of B at the middle of t^{th} iteration (before update step $\tilde{x}^{(t+1)} = A(y^{(t)})^{*2}$) and $B^{(t,2)}$ is the conditional distribution at the end of t^{th} iteration (after update step $\tilde{x}^{(t+1)} = A(y^{(t)})^{*2}$). By analyzing the above conditional distributions, we can characterize the left independent randomness in B .

3.1.1 Conditional Distributions

In order to characterize the conditional distribution of B under evolution of $x^{(t)}$ and $y^{(t)}$ in (15), we exploit the following basic fact (see (Bayati and Montanari, 2010) for proof).

Lemma 7 (Conditional distribution of Gaussian matrices under a linear constraint). *Consider random matrix D with i.i.d. Gaussian entries $D_{ij} \sim \mathcal{N}(0, \sigma^2)$. Conditioned on $u = Dv$ with known vectors u and v , the matrix D is distributed as*

$$D|\{u = Dv\} \stackrel{(d)}{=} \frac{1}{\|v\|^2} uv^\top + \tilde{D}P_{\perp v},$$

where random matrix \tilde{D} is an independent copy of D with i.i.d. Gaussian entries $\tilde{D}_{ij} \sim \mathcal{N}(0, \sigma^2)$, and $P_{\perp v}$ is the projection operator on to the subspace orthogonal to v .

We refer to $\tilde{D}P_{\perp v}$ as the *residual* random matrix since it represents the remaining *randomness* left after conditioning. It is a random matrix whose rows are independent random vectors that are orthogonal to v , and the variance in each direction orthogonal to v is equal to σ^2 .

The above Lemma can be exploited to characterize the conditional distribution of B introduced in (15). However, a naive direct application using the constraint $Y^{[t]} = A^\top X^{[t]}$ is not transparent for analysis. The reason is the evolution of $x^{(t)}$ and $y^{(t)}$ are themselves governed by the conditional distribution of B given previous iterations. Therefore, we need the following recursive version of Lemma 7.

Corollary 1 (Iterative conditioning). *Consider random matrix D with i.i.d. Gaussian entries $D_{ij} \sim \mathcal{N}(0, \sigma^2)$, and let $F \stackrel{(d)}{=} P_{\perp C} D P_{\perp R}$ be the random Gaussian matrix whose columns are orthogonal to space C and rows are orthogonal to space R . Conditioned on the linear constraint $u = Dv$, where⁹ $u \in C^\perp$, the matrix F is distributed as*

$$F|\{u = Dv\} \stackrel{(d)}{=} \frac{1}{\|(P_{\perp R}v)\|^2} u(P_{\perp R}v)^\top + P_{\perp C} \tilde{D} P_{\perp \{R, v\}},$$

where random matrix \tilde{D} is an independent copy of D with i.i.d. Gaussian entries $\tilde{D}_{ij} \sim \mathcal{N}(0, \sigma^2)$.

Thus, the *residual* random matrix $P_{\perp C} \tilde{D} P_{\perp \{R, v\}}$ is a random Gaussian matrix whose columns are orthogonal to C and rows are orthogonal to $\text{span}\{R, v\}$. The variance in any remaining dimension is equal to σ^2 .

3.1.2 Form of Iterative Updates

Now we exploit the conditional distribution arguments proposed in the previous section to characterize the conditional distribution of B given the update variables x and y up to the current iteration; recall (15) where $B^{(t,1)}$ is the conditional distribution of B at the middle of t^{th} iteration and $B^{(t,2)}$ at the end of t^{th} iteration. Before that, we need to introduce some more intermediate variables.

Intermediate variables: We separate the first entry of y and $(y)^{*2}$ from the rest, i.e., we have

$$y_1^{(t)} = a_1^\top x^{(t)}, \quad y_{-1}^{(t)} = B^\top x^{(t)} \sim (B^{(t-1,2)})^\top x^{(t)},$$

where $y_{-1}^{(t)} \in \mathbb{R}^{k-1}$ denotes $y^{(t)} \in \mathbb{R}^k$ with the first entry removed. The update formula for $\tilde{x}^{(t+1)}$ can be also decomposed as

$$\tilde{x}^{(t+1)} = (y_1^{(t)})^2 a_1 + B w^{(t)} \sim (y_1^{(t)})^2 a_1 + B^{(t,1)} w^{(t)},$$

where

$$w^{(t)} := (y_{-1}^{(t)})^{*2} \in \mathbb{R}^{k-1},$$

is the new intermediate variable in the power iterations. Let $B_{\text{res.}}^{(t,1)}$ and $B_{\text{res.}}^{(t,2)}$ denote the *residual* random matrices corresponding to $B^{(t,1)}$ and $B^{(t,2)}$ respectively, and

$$u^{(t+1)} := B_{\text{res.}}^{(t,1)} w^{(t)}, \quad v^{(t)} := (B_{\text{res.}}^{(t-1,2)})^\top x^{(t)},$$

⁹We need that $u \in C^\perp$, otherwise the event $u = Dv$ is impossible.

where $u^{(t)} \in \mathbb{R}^d$ and $v^{(t)} \in \mathbb{R}^{k-1}$ are respectively the part of $x^{(t)}$ and $y_{-1}^{(t)}$ representing the residual randomness after conditioning on the previous iterations. We also summarize all variables and notations in Table 1 in the Appendix which can be used as a reference throughout the paper.

Finally we make the following observations.

Lemma 8 (Form of iterative updates). *The conditional distribution of B at the middle of t^{th} iteration denoted by $B^{(t,1)}$ satisfies*

$$B^{(t,1)} \stackrel{(d)}{=} \sum_{i \in [t-1]} \frac{u^{(i+1)}(P_{\perp_W^{[i-1]}} w^{(i)})^\top}{\|P_{\perp_W^{[i-1]}} w^{(i)}\|^2} + \sum_{i \in [t]} \frac{P_{\perp_X^{[i-1]}} x^{(i)}(v^{(i)})^\top}{\|P_{\perp_X^{[i-1]}} x^{(i)}\|^2} + B_{\text{res.}}^{(t,1)}, \quad (16)$$

$$B_{\text{res.}}^{(t,1)} \stackrel{(d)}{=} P_{\perp_{X^{[t]}}} \tilde{B} P_{\perp_{W^{[t-1]}}}, \quad (17)$$

where random matrix \tilde{B} is an independent copy of B with i.i.d. Gaussian entries $\tilde{B}_{ij} \sim \mathcal{N}(0, \frac{1}{d})$. Similarly, the conditional distribution of B at the end of t^{th} iteration denoted by $B^{(t,2)}$ satisfies

$$B^{(t,2)} \stackrel{(d)}{=} \sum_{i \in [t]} \left(\frac{u^{(i+1)}(P_{\perp_W^{[i-1]}} w^{(i)})^\top}{\|P_{\perp_W^{[i-1]}} w^{(i)}\|^2} + \frac{P_{\perp_X^{[i-1]}} x^{(i)}(v^{(i)})^\top}{\|P_{\perp_X^{[i-1]}} x^{(i)}\|^2} \right) + B_{\text{res.}}^{(t,2)}, \quad (18)$$

$$B_{\text{res.}}^{(t,2)} \stackrel{(d)}{=} P_{\perp_{X^{[t]}}} B' P_{\perp_{W^{[t]}}}, \quad (19)$$

where random matrix B' is an independent copy of B with i.i.d. Gaussian entries $B'_{ij} \sim \mathcal{N}(0, \frac{1}{d})$.

The lemma can be directly proved by applying the iterative conditioning argument in Corollary 1. See the detailed proof in the appendix.

3.1.3 Analysis of Iterative Updates

Lemma 8 characterizes the conditional distribution of B given the update variables x and y up to the current iteration; see (15) for the definition of conditional forms of B denoted by $B^{(t,1)}$ and $B^{(t,2)}$. Intuitively, when the number of iterations $t \ll d$, then the residual independent randomness left in $B^{(t,1)}$ and $B^{(t,2)}$ (respectively denoted by $B_{\text{res.}}^{(t,1)}$ and $B_{\text{res.}}^{(t,2)}$) characterized in Lemma 8 is large enough and we can exploit it to obtain tighter concentration bounds throughout the analysis of the iterations.

Note that the goal is to show that under $t \ll d$, the iterations $x^{(t)}$ converge to the true component with constant error, i.e., $|\langle x^{(t)}, a_1 \rangle| \geq 1 - \gamma$ for some constant $\gamma > 0$. If this already holds before iteration t we are done, and if it does not hold, next iteration is analyzed to finally achieve the goal. This analysis is done via *induction argument*. During the iterations, we maintain several invariants to analyze the dynamics of power update. The goal is to ensure progress in each iteration as in (20).

Induction hypothesis: The following are assumed at the beginning of the iteration t as induction hypothesis; see Figure 2 for the scope of inductive step.

1. Length of Projection on x :

$$\delta_t \leq \|P_{\perp_{X^{[t-1]}}} x^{(t)}\| \leq 1,$$

where δ_t is of order $1/\text{polylog } d$, and the value of δ_t only depends on t and $\log d$.

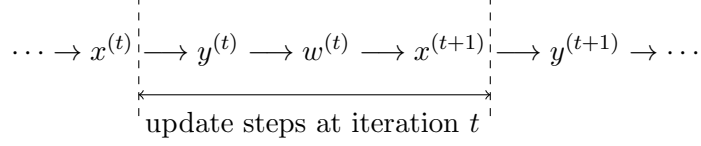


Figure 2: Flow of the power update algorithm stating intermediate steps. Iteration t for which the inductive step should be argued is also indicated.

2. Length of Projection on w :

$$\begin{aligned} \delta'_{t-1} \frac{\sqrt{k}}{d} &\leq \|P_{\perp_{W^{[t-2]}}} w^{(t-1)}\| \leq \Delta'_{t-1} \frac{\sqrt{k}}{d}, \\ \|P_{\perp_{W^{[t-2]}}} w^{(t-1)}\|_{\infty} &\leq \Delta'_{t-1} \frac{1}{d}, \end{aligned}$$

where δ'_t is of order $1/\text{polylog } d$ and Δ'_t is of order $\text{polylog } d$. Both δ'_t and Δ'_t only depend on t and $\log d$.

3. Progress:¹⁰

$$\begin{aligned} |\langle a_1, x^{(t)} \rangle| &\in [\delta_t^*, \Delta_t^*] d^{\beta 2^{t-1}} \frac{\sqrt{k}}{d}, \\ \langle a_1, P_{\perp_{X^{[t-1]}}} x^{(t)} \rangle &\leq \Delta_t^* d^{\beta 2^{t-1}} \frac{\sqrt{k}}{d}. \end{aligned} \tag{20}$$

4. Norm of u, v :

$$\begin{aligned} \frac{\delta_{t-1}}{2} \sqrt{\frac{k}{d}} &\leq \|v^{(t-1)}\| \leq 2\sqrt{\frac{k}{d}}, \\ \frac{\delta'_{t-1}}{2} \frac{\sqrt{k}}{d} &\leq \|u^{(t)}\| \leq 2\Delta'_{t-1} \frac{\sqrt{k}}{d}. \end{aligned}$$

The analysis for basis of induction and inductive step are provided in Appendix A.

3.2 Effect of noise in Theorem 1

Given rank- k random tensor T in (12), and a starting point $x^{(1)}$, our analysis in the noiseless setting shows that the tensor power iteration in (2) outputs a vector which will be close to a_j if $x^{(1)}$ has a large enough correlation with a_j .

Now suppose we are given noisy tensor $\hat{T} = T + E$ where E has some small norm. In this case where the noise is also present, we get a sequence $\hat{x}^{(t)} = x^{(t)} + \xi^{(t)}$ where $x^{(t)}$ is the component not incorporating any noise (as in previous section¹¹), while $\xi^{(t)}$ represents the contribution of noise tensor E in the power iteration; see (21) below. We prove that $\xi^{(t)}$ is a very small noise that does not change our calculations stated in the following lemma.

¹⁰Note that although the bounds on $y_{-1}^{(t)}$ are argued at iteration t , the bound on the first entry of $y^{(t)}$ denoted by $y_1^{(t)} = \langle a_1, x^{(t)} \rangle$ is assumed here in the induction hypothesis at the end of iteration $t - 1$.

¹¹Note that there is a subtle difference between notation $x^{(t)}$ in the noiseless and noisy settings. In the noiseless setting, this vector is normalized, while in the noisy setting the whole vector $\hat{x}^{(t)} = x^{(t)} + \xi^{(t)}$ is normalized.

Lemma 9 (Bounding norm of error). *Suppose the spectral norm of the error tensor E is bounded as*

$$\|E\| \leq \epsilon\sqrt{k}/d, \quad \text{where } \epsilon < o(\sqrt{k}/d).$$

Then the noise vector $\xi^{(t)}$ at iteration t satisfies the ℓ_2 norm bound

$$\|\xi^{(t)}\| \leq \tilde{O}(d^{\beta 2^{t-1}} \epsilon).$$

Note that when t is the first number such that $d^{\beta 2^{t-1}} \geq d/\sqrt{k}$, we have $\|\xi^{(t)}\| = o(1)$.

Notice that since when $d^{\beta 2^{t-1}} \geq d/\sqrt{k}$, the main induction is already over and we know $x^{(t)}$ is constant close to the true component, and thus, the noise is always small.

Proof idea: We now provide an overview of ideas for proving the above lemma; see Appendix C for the complete proof which is based on an induction argument. We first write the following recursion expanding the contribution of main signal and noise terms in the tensor power iteration as

$$\begin{aligned} x^{(t+1)} + \xi^{(t+1)} &= \text{Norm} \left(\hat{T}(x^{(t)} + \xi^{(t)}, x^{(t)} + \xi^{(t)}, I) \right) \\ &= \text{Norm} \left(T(x^{(t)}, x^{(t)}, I) + 2T(x^{(t)}, \xi^{(t)}, I) + T(\xi^{(t)}, \xi^{(t)}, I) + E(\hat{x}^{(t)}, \hat{x}^{(t)}, I) \right), \end{aligned} \quad (21)$$

where for vector v , we have $\text{Norm}(v) := v/\|v\|$, i.e., it normalizes the vector. The first term is the desired main signal and should have the largest norm, and the rest of the terms are the noise terms. The third term is of order $\|\xi^{(t)}\|^2$, and hence, it should be fine whenever we choose $\|E\|$ to be small enough. The last term is $O(\|E\|)$ and is the same for all iterations so that is also fine. The problematic term is the second term, whose norm if we bound naively is $2\|\xi^{(t)}\|$. However the normalization factor also contributes a factor of roughly d/\sqrt{k} , and thus, this term grows exponentially; it is still fine if we just do a constant number of iterations, but the exponent will depend on the number of iterations.

In order to solve this problem, and make sure that the amount of noise we can tolerate is independent of the number of iterations, we need a better way to bound the noise term $\xi^{(t)}$. The main problem here is we bound the norm of $\|T(x^{(t)}, \xi^{(t)}, I)\|$ by $\|T\|\|\xi^{(t)}\| \leq O(\xi^{(t)})$, by doing this we ignored the fact that $x^{(t)}$ is uncorrelated with the components in T . In order to get a tighter bound, we introduce another norm $\|\cdot\|_*$. Intuitively, the norm $\|\cdot\|_*$ captures the fact that x does not have a high correlation with the components (except for the first component that x will converge to), and gives a better bound. In particular we have $\|T(x^{(t)}, \xi^{(t)}, I)\| \approx \frac{\sqrt{k}}{d} \|\xi^{(t)}\|_2$. Therefore, the normalization factor is compensated by the additional term $\frac{\sqrt{k}}{d}$. More concretely, this norm is defined as follows.

Definition 1 (Norm $\|\cdot\|_*$). *Given a matrix $A = [a_1 \ a_2 \ \dots \ a_k] \in \mathbb{R}^{d \times k}$, for any vector $u \in \mathbb{R}^d$, the norm $\|u\|_{A^*}$ is defined as*

$$\|u\|_{A^*} = \max_{i \in [k]} |\langle a_i, u \rangle|.$$

This norm satisfies a property shown in Lemma 19 which enables us to argue that $\xi^{(t)}$ is small enough as stated in Lemma 9.

4 Conclusion

In this paper, we provide a novel analysis for the dynamics of third order tensor power iteration showing convergence guarantees to vectors having constant correlation with the tensor component. This enables us to prove unsupervised learning of latent variable models in the challenging over-complete regime where the hidden dimensionality is larger than the observed dimension. The main technical observation is that under random Gaussian tensor components and small number of iterations, the residual randomness in the components (which are involved in the iterative steps) are sufficiently large. This enables us to show progress in the next iteration of the update step. As future work, it is very interesting to generalize this analysis to higher order tensor power iteration, and more generally to other kinds of iterative updates.

Acknowledgements

A. Anandkumar is supported in part by Microsoft Faculty Fellowship, NSF Career award CCF-1254106, NSF Award CCF-1219234, ARO YIP Award W911NF-13-1-0084 and ONR Award N00014-14-1-0665. M. Janzamin is supported by NSF Award CCF-1219234.

Appendix

Table 1: Table of parameters and variables. Superscript (t) denotes the variable at t -th iteration.

Variable	Space	Description	Recursion formula
A	$\mathbb{R}^{d \times k}$	mapping matrix in update formula (14)	n.a.
$x^{(t)}$	\mathbb{R}^d	update variable in (14)	$x^{(t+1)} := \frac{A(y^{(t)})^{*2}}{\ A(y^{(t)})^{*2}\ }$
$y^{(t)}$	\mathbb{R}^k	intermediate variable in update formula (14)	$y^{(t)} := A^\top x^{(t)}$
$\tilde{x}^{(t)}$	\mathbb{R}^d	unnormalized version of $x^{(t)}$	$\tilde{x}^{(t+1)} := A(y^{(t)})^{*2}$
$\hat{x}^{(t)}$	\mathbb{R}^d	noisy version of $x^{(t)}$	$\hat{x}^{(t)} = x^{(t)} + \xi^{(t)}$; see (21)
$\xi^{(t)}$	\mathbb{R}^d	Contribution of noise in tensor power update given noisy tensor $\hat{T} = T + E$	$\hat{x}^{(t)} = x^{(t)} + \xi^{(t)}$; see (21)
B	$\mathbb{R}^{d \times (k-1)}$	matrix $A := [a_1 \ a_2 \ \cdots \ a_k]$ with first column removed, i.e., $B := [a_2 \ a_3 \ \cdots \ a_k]$. Note that the first column a_1 is the desired one to recover.	n.a.
$B^{(t,1)}$	$\mathbb{R}^{d \times (k-1)}$	conditional distribution of B given previous iterations at the middle of t^{th} iteration (before update step $\tilde{x}^{(t+1)} = A(y^{(t)})^{*2}$)	$B^{(t,1)} \stackrel{(d)}{=} B \{X^{[t]}, Y^{[t]}\}$
$B^{(t,2)}$	$\mathbb{R}^{d \times (k-1)}$	conditional distribution of B given previous iterations at the end of t^{th} iteration (after update step $\tilde{x}^{(t+1)} = A(y^{(t)})^{*2}$)	$B^{(t,2)} \stackrel{(d)}{=} B \{X^{[t+1]}, Y^{[t]}\}$
$B_{\text{res.}}^{(t,1)}$	$\mathbb{R}^{d \times (k-1)}$	residual independent randomness left in $B^{(t,1)}$; see Lemma 8.	see equation (17)
$B_{\text{res.}}^{(t,2)}$	$\mathbb{R}^{d \times (k-1)}$	residual independent randomness left in $B^{(t,2)}$; see Lemma 8.	see equation (19)
$w^{(t)}$	\mathbb{R}^{k-1}	intermediate variable in update formula (14)	$w^{(t)} := (y_{-1}^{(t)})^{*2}$
$u^{(t)}$	\mathbb{R}^d	part of $x^{(t)}$ representing the left independent randomness	$u^{(t+1)} := B_{\text{res.}}^{(t,1)} w^{(t)}$
$v^{(t)}$	\mathbb{R}^{k-1}	part of $y_{-1}^{(t)}$ representing the left independent randomness	$v^{(t)} := (B_{\text{res.}}^{(t-1,2)})^\top x^{(t)}$

Proof of Lemma 8: Recall that we have updates of the form

$$\tilde{x}^{(t+1)} = A(y^{(t)})^{*2}, \quad w^{(t)} := (y_{-1}^{(t)})^{*2}, \quad y^{(t)} = A^\top x^{(t)}.$$

Let

$$X^{[t]\setminus 1} := \left[x^{(2)} \mid \dots \mid x^{(t)} \right],$$

and let the rows of $Y^{[t]}$ are partitioned as the first and the rest of rows as

$$Y^{[t]} = \left[Y_1^{[t]\top} \mid Y_{-1}^{[t]\top} \right]^\top.$$

We now make the following simple observations

$$\begin{aligned} B^{(t,1)} &\stackrel{(d)}{=} B \{ Y^{[t]} = A^\top X^{[t]}, \tilde{X}^{[t]\setminus 1} = A(Y^{[t-1]})^{*2} \} \\ &\stackrel{(d)}{=} B \{ Y_{-1}^{[t]} = B^\top X^{[t]}, \tilde{X}^{[t]\setminus 1} = a_1(Y_1^{[t-1]})^{*2} + BW^{[t-1]} \} \\ &\stackrel{(d)}{=} B \{ v^{(1)} = B^\top x^{(1)}, \dots, v^{(t)} = (B_{\text{res.}}^{(t-1,2)})^\top x^{(t)}, \\ &\quad u^{(2)} = B_{\text{res.}}^{(1,1)} w^{(1)}, \dots, u^{(t)} = B_{\text{res.}}^{(t-1,1)} w^{(t-1)} \}, \end{aligned}$$

where the second equivalence comes from the fact that B is matrix A with first column removed. Now applying Corollary 1, we have the result. The distribution of $B^{(t,2)}$ follow similarly. \square

A Analysis of Induction Argument

In this section, we analyze the basis of induction and inductive step for the induction argument proposed in Section 3.1.3 for the proof of Lemma 5.

A.1 Basis of induction

We first show that the hypothesis holds for initialization vector $x^{(1)}$ as the basis of induction.

Claim 1 (Basis of induction). *The induction hypothesis is true for $t = 1$.*

Proof: Notice that induction hypothesis for $t = 1$ only involves the bounds on $\|x^{(1)}\|$ and $\langle a_1, x^{(1)} \rangle$ as in Hypotheses 1 and 3, respectively. These bounds are directly argued by the correlation assumption on the initial vector $x^{(1)}$ stated in (13) where $\delta_1 = \delta_1^* = \Delta_1^* = 1$. \square

A.2 Inductive step

Assuming the induction hypothesis holds for all the values till the end of iteration $t - 1$ (stated in Section 3.1.3), we analyze the t -th iteration of the algorithm, and prove that induction hypothesis also holds for the values at the end of iteration t . See Figure 2 where the scope of iteration t and the flow of the algorithm is shown. In the rest of this section, we pursue the flow of the algorithm at iteration t starting from computing $y^{(t)}$ and ending up with computing $x^{(t+1)}$ to prove the desired induction hypothesis at iteration t .

Hypothesis 4

We start by showing that the induction Hypothesis 4 holds at iteration t using the induction Hypotheses 1 and 2 in the previous iteration.

Claim 2. *We have*

$$\begin{aligned}\frac{\delta_t}{2}\sqrt{\frac{k}{d}} &\leq \|v^{(t)}\| \leq 2\sqrt{\frac{k}{d}}, \\ \frac{\delta'_t}{2}\frac{\sqrt{k}}{d} &\leq \|u^{(t+1)}\| \leq 2\Delta'_t\frac{\sqrt{k}}{d}.\end{aligned}$$

Proof: Recall that $v^{(t)} := (B_{\text{res.}}^{(t-1,2)})^\top x^{(t)}$, and by applying the form of $B_{\text{res.}}^{(t-1,2)}$ in (19), we have

$$v^{(t)} \stackrel{(d)}{=} P_{\perp_{W^{[t-1]}}} B'^\top P_{\perp_{X^{[t-1]}}} x^{(t)}. \quad (22)$$

Since random matrix $B' \in \mathbb{R}^{d \times (k-1)}$ is an independent copy of B with i.i.d. Gaussian entries $B'_{ij} \sim \mathcal{N}(0, \frac{1}{d})$, we know $v^{(t)}$ is a random Gaussian vector in the subspace orthogonal to $W^{[t-1]}$. On the other hand, for any vector $z \in \mathbb{R}^d$, we have

$$\mathbb{E} \left[\|B'^\top z\|^2 \right] = z^\top \mathbb{E} \left[B' B'^\top \right] z = \frac{k-1}{d} \|z\|^2,$$

where $\mathbb{E} [B' B'^\top] = \frac{k-1}{d} I$ is exploited. Let $z = P_{\perp_{X^{[t-1]}}} x^{(t)}$. Then, by applying the above equality to the expansion of $v^{(t)}$ in (22), we have

$$\mathbb{E} \left[\|v^{(t)}\|^2 \right] = \frac{k-t}{k-1} \cdot \frac{k-1}{d} \cdot \|P_{\perp_{X^{[t-1]}}} x^{(t)}\|^2 = \frac{k-t}{d} \cdot \|P_{\perp_{X^{[t-1]}}} x^{(t)}\|^2 \in \left[\delta_t^2 \frac{k}{d} \left(1 - \frac{t}{k}\right), \frac{k}{d} \right],$$

where $\dim(W^{[t-1]}) = t-1$ is also used in the first step, and the last step is concluded from Hypothesis 1. Finally, by concentration property of random Gaussian vectors, when $t \ll d$ we have with high probability

$$\|v^{(t)}\| \in \left[\frac{\delta_t}{2}\sqrt{\frac{k}{d}}, 2\sqrt{\frac{k}{d}} \right].$$

Similarly, for $u^{(t+1)} := B_{\text{res.}}^{(t,1)} w^{(t)}$, and by applying the form of $B_{\text{res.}}^{(t,1)}$ in (17), we have

$$u^{(t+1)} \stackrel{(d)}{=} P_{\perp_{X^{[t]}}} \tilde{B} P_{\perp_{W^{[t-1]}}} w^{(t)}. \quad (23)$$

Since random matrix $\tilde{B} \in \mathbb{R}^{d \times (k-1)}$ is an independent copy of B with i.i.d. Gaussian entries $\tilde{B}_{ij} \sim \mathcal{N}(0, \frac{1}{d})$, we know $u^{(t+1)}$ is a random Gaussian vector in the subspace orthogonal to $X^{[t]}$. On the other hand, for any vector $z \in \mathbb{R}^{k-1}$, we have

$$\mathbb{E} \left[\|\tilde{B} z\|^2 \right] = z^\top \mathbb{E} \left[\tilde{B}^\top \tilde{B} \right] z = \|z\|^2,$$

where $\mathbb{E} [\tilde{B}^\top \tilde{B}] = I$ is exploited. Let $z = P_{\perp_{W^{[t-1]}}} w^{(t)}$. Then, by applying the above equality to the expansion of $u^{(t+1)}$ in (23), we have

$$\mathbb{E} \left[\|u^{(t+1)}\|^2 \right] = \frac{d-t}{d} \cdot \|P_{\perp_{W^{[t-1]}}} w^{(t)}\|^2 \in \left[(\delta'_t)^2 \frac{k}{d^2} \left(1 - \frac{t}{d}\right), (\Delta'_t)^2 \frac{k}{d^2} \right],$$

where $\dim(X^{[t]}) = t$ is also used in the first step, and the last step is concluded from Hypothesis 2. Finally, by concentration property of random Gaussian vectors, when $t \ll d$ we have with high probability

$$\|u^{(t+1)}\| \in \left[\frac{\delta'_t \sqrt{k}}{2d}, 2\Delta'_t \frac{\sqrt{k}}{d} \right].$$

□

Hypothesis 2

Computing $y^{(t)}$: In the first step of iteration t , the algorithm computes $y^{(t)}$. By induction Hypothesis 3, we know $|y_1^{(t)}| = \tilde{\Theta}(d^{\beta 2^{t-1}} \sqrt{k}/d)$. The other coordinates of $y^{(t)} := A^\top x^{(t)}$ are $y_{-1}^{(t)} = B^\top x^{(t)}$ which conditioning on the previous iterations are equivalent (in distribution) to

$$\begin{aligned} y_{-1}^{(t)} &\stackrel{(d)}{=} \left(B^{(t-1,2)} \right)^\top x^{(t)} \\ &= \left(\sum_{i \in [t-1]} \left(\frac{u^{(i+1)} (P_{\perp_W^{[i-1]}} w^{(i)})^\top}{\|P_{\perp_W^{[i-1]}} w^{(i)}\|^2} + \frac{P_{\perp_X^{[i-1]}} x^{(i)} (v^{(i)})^\top}{\|P_{\perp_X^{[i-1]}} x^{(i)}\|^2} \right) + B_{\text{res.}}^{(t-1,2)} \right)^\top x^{(t)} \\ &= \sum_{i \in [t-1]} \left(\tilde{\Theta} \left(\frac{d^2}{k} \right) P_{\perp_W^{[i-1]}} w^{(i)} \langle u^{(i+1)}, x^{(t)} \rangle + \tilde{\Theta}(1) v^{(i)} \langle P_{\perp_X^{[i-1]}} x^{(i)}, x^{(t)} \rangle \right) + v^{(t)}, \end{aligned} \quad (24)$$

where form of $B^{(t-1,2)}$ in (18) is used in the second equality. The bounds on the norms come from Hypotheses 1 and 2. The last term is by definition $v^{(t)} := (B_{\text{res.}}^{(t-1,2)})^\top x^{(t)}$. Note that differences in polylog factors in the (upper and lower) bounds in Hypotheses 1 and 2 are represented by notation $\tilde{\Theta}(\cdot)$.

We will establish subsequently that if $k > d$, the terms involving $v^{(i)}$'s in the above expansion dominate, and the terms involving $P_{\perp_W^{[i-1]}} w^{(i)}$'s have norm of a smaller order; see Claim 3.

Computing $w^{(t)}$: In the next step of the algorithm at iteration t , $w^{(t)}$ is computed for which we now argue if the induction hypothesis holds up to iteration t , both lower and upper bounds at iteration t as $\|P_{\perp_W^{[t-1]}} w^{(t)}\| \in [\delta'_t, \Delta'_t] \frac{\sqrt{k}}{d}$ (see induction Hypothesis 2) also hold.

Lower bound: For the lower bound, intuitively the *fresh* random vector $v^{(t)}$ should bring enough randomness into $w^{(t)}$. We formulate that in the following lemma.

Lemma 10. *Suppose R and R' are two subspaces in \mathbb{R}^k with dimension at most $t \leq \frac{k}{16(\log k)^2}$. Let $p \in \mathbb{R}^k$ be an arbitrary vector, $z \in \mathbb{R}^k$ be a uniformly random Gaussian vector in the space orthogonal to R , and finally $w = (p + z) * (p + z)$. Then with high probability, we have*

$$\|P_{\perp_{R'}} w\| \geq \frac{\mathbb{E}[\|z\|^2]}{40\sqrt{k}}.$$

Recall that $w^{(t)} := y_{-1}^{(t)} * y_{-1}^{(t)}$, and $y_{-1}^{(t)}$ is expanded in (24) as sum of an arbitrary vector and a random Gaussian vector. Applying above lemma with $R = R' = \text{span}(W^{[t-1]})$, we have with high probability

$$\|P_{\perp_{W^{[t-1]}}} w^{(t)}\| \geq \frac{\mathbb{E}[\|v^{(t)}\|^2]}{40\sqrt{k}} \geq \frac{\delta_t^2}{160} \sqrt{k}/d,$$

where Hypothesis 4 gives lower bound $\|v^{(t)}\| \geq \delta_t/2\sqrt{k/d}$ (used in the second inequality). By choosing $\delta_t' = \delta_t^2/160$ the lower bound in Hypothesis 2 is proved.

Upper bound: In order to prove the upper bounds in Hypothesis 2, we follow the sequence of arguments below:

$$\text{Claim 3: } \|y_{-1}^{(t)}\|_{\infty} \xrightarrow{(\cdot)^2} \|w^{(t)}\|_{\infty} \xrightarrow{\text{Lemma 11}} \|P_{\perp_{W^{[t-1]}}} w^{(t)}\|_{\infty} \Rightarrow \|P_{\perp_{W^{[t-1]}}} w^{(t)}\|$$

First we prove a bound on the infinity norm of $y_{-1}^{(t)}$:

Claim 3 (Upper bound on $\|y_{-1}^{(t)}\|_{\infty}$). *We have*

$$\|y_{-1}^{(t)}\|_{\infty} \leq \frac{t \log d}{\delta_t \sqrt{d}} + (t-1) \left(\frac{\Delta'_{t-1}}{\delta'_{t-1}} \right)^2 \frac{1}{\sqrt{k}} = \tilde{O} \left(\frac{1}{\sqrt{d}} \right).$$

Proof: We exploit the induction hypothesis to bound the ℓ_{∞} norm of all the terms in the expansion of $y_{-1}^{(t)}$ in (24).

For the terms involving $v^{(i)}$, since they are random Gaussian vectors with expected square norm at most k/d , by Lemma 17 we know $\|v^{(i)}\|_{\infty} \leq \frac{\log d}{\sqrt{d}}$ with high probability. In addition, for $v^{(i)}$, $i < t$, the coefficient is bounded as

$$\frac{\langle P_{\perp_{X^{[i-1]}}} x^{(i)}, x^{(t)} \rangle}{\|P_{\perp_{X^{[i-1]}}} x^{(i)}\|^2} \leq \frac{1}{\|P_{\perp_{X^{[i-1]}}} x^{(i)}\|} \leq \frac{1}{\delta_i}, \quad (25)$$

where the last step uses Hypothesis 1. Therefore, the total contribution from terms involving $v^{(i)}$ in $\|y_{-1}^{(t)}\|_{\infty}$ is bounded by $\frac{t \log d}{\delta_t \sqrt{d}}$.

For the terms involving $P_{\perp_{W^{[i-1]}}} w^{(i)}$, $i \in [t-1]$, we have from Hypothesis 2 that the ℓ_{∞} norm is bounded as $\|P_{\perp_{W^{[i-1]}}} w^{(i)}\|_{\infty} \leq \Delta_i' \frac{1}{d}$. In addition, the corresponding coefficient is bounded by

$$\frac{\langle u^{(i+1)}, x^{(t)} \rangle}{\|P_{\perp_{W^{[i-1]}}} w^{(i)}\|^2} \leq \frac{\|u^{(i+1)}\| \cdot \|x^{(t)}\|}{\|P_{\perp_{W^{[i-1]}}} w^{(i)}\|^2} \leq \frac{2\Delta_i'}{\delta_i^2} \frac{d}{\sqrt{k}}. \quad (26)$$

Again bounds in Hypotheses 2 and 4 are exploited in the last inequality. Hence, the total contribution from terms involving $P_{\perp_{W^{[i-1]}}} w^{(i)}$, $i \in [t-1]$ in $\|y_{-1}^{(t)}\|_{\infty}$ is bounded by $(t-1) \left(\frac{\Delta'_{t-1}}{\delta'_{t-1}} \right)^2 \frac{1}{\sqrt{k}}$.

Combining the above bounds finishes the proof. \square

Since $w^{(t)} := y_{-1}^{(t)} * y_{-1}^{(t)}$, the above claim immediately implies that

$$\|w^{(t)}\|_{\infty} \leq \tilde{O} \left(\frac{1}{d} \right). \quad (27)$$

Now we have the ℓ_∞ norm on w , however we need to bound the ℓ_∞ norm of the projected vector $P_{\perp_{W^{[t-1]}}} w^{(t)}$. Intuitively this is clear as the vectors in the space $W^{[t-1]}$ all have small ℓ_∞ as guaranteed by induction hypothesis. We formalize this intuition using the following lemma.

Lemma 11. *Suppose R is a subspace in \mathbb{R}^k of dimension t' , such that there is a basis $\{r_1, \dots, r_{t'}\}$ with $\|r_i\|_\infty \leq \frac{\Delta}{\sqrt{k}}$ and $\|r_i\| = 1$. Let $p \in \mathbb{R}^k$ be an arbitrary vector, then*

$$\|P_{\perp_R} p\|_\infty \leq \|p\|_\infty + \|p\| \Delta \frac{\sqrt{t'}}{\sqrt{k}}.$$

Let $R = \text{span}(W^{[t-1]})$. Then the vectors $P_{\perp_{W^{[i-1]}}} w^{(i)} / \|P_{\perp_{W^{[i-1]}}} w^{(i)}\|$, $i \in [t-1]$ form a basis for subspace R , and we know from Hypothesis 2 that the ℓ_∞ norm of each of these basis vectors is bounded by $\frac{\Delta}{\sqrt{k}}$ for $\Delta := \frac{\Delta'_{t-1}}{\delta'_{t-1}}$ which is of order polylog d . Applying above lemma, we have

$$\|P_{\perp_{W^{[t-1]}}} w^{(t)}\|_\infty \leq \|w^{(t)}\|_\infty (1 + \Delta \sqrt{t-1}) \leq \frac{\Delta'_t}{d},$$

where the last inequality uses bound (27), and appropriate choosing for Δ'_t which is of order polylog d and only depends on t and $\log d$. This concludes the upper bound on the ℓ_∞ norm in Hypothesis 2. The upper bound on the ℓ_2 norm is also immediately argued using this ℓ_∞ norm bound where an additional \sqrt{k} factor shows up.

Hypothesis 1

Computing $x^{(t+1)}$: In the next step of iteration t , the algorithm computes $x^{(t+1)}$. Conditioning on the previous iterations, the unnormalized version $\tilde{x}^{(t+1)}$ is equivalent (in distribution) to

$$\begin{aligned} \tilde{x}^{(t+1)} &\stackrel{(d)}{=} B^{(t,1)} w^{(t)} + (y_1^{(t)})^2 a_1 \\ &= \sum_{i \in [t-1]} \frac{u^{(i+1)} (P_{\perp_{W^{[i-1]}}} w^{(i)})^\top}{\|P_{\perp_{W^{[i-1]}}} w^{(i)}\|^2} w^{(t)} + \sum_{i \in [t]} \frac{P_{\perp_{X^{[i-1]}}} x^{(i)} (v^{(i)})^\top}{\|P_{\perp_{X^{[i-1]}}} x^{(i)}\|^2} w^{(t)} + B_{\text{res.}}^{(t,1)} w^{(t)} + (y_1^{(t)})^2 a_1 \\ &= \sum_{i \in [t-1]} \tilde{\Theta} \left(\frac{d^2}{k} \right) u^{(i+1)} \langle P_{\perp_{W^{[i-1]}}} w^{(i)}, w^{(t)} \rangle + \sum_{i \in [t]} \tilde{\Theta}(1) P_{\perp_{X^{[i-1]}}} x^{(i)} \langle v^{(i)}, w^{(t)} \rangle + u^{(t+1)} + (y_1^{(t)})^2 a_1, \end{aligned} \tag{28}$$

where form of $B^{(t,1)}$ in (16) is used in the second equality. The bounds on the norms come from Hypotheses 1 and 2. The last term is the definition of $u^{(t+1)} := B_{\text{res.}}^{(t,1)} w^{(t)}$. Note that differences in polylog factors in the (upper and lower) bounds in Hypotheses 1 and 2 are represented by notation $\tilde{\Theta}(\cdot)$.

The goal is to prove Hypothesis 1 holds at t -th iteration (which is to show the desired lower and upper bounds on $\|P_{\perp_{X^{[t]}}} x^{(t+1)}\|$) assuming induction hypothesis holds for earlier iterations. Given the normalization $x^{(t+1)} := \tilde{x}^{(t+1)} / \|\tilde{x}^{(t+1)}\|$ in each iteration, we have

$$\|P_{\perp_{X^{[t]}}} x^{(t+1)}\| = \frac{1}{\|\tilde{x}^{(t+1)}\|} \|P_{\perp_{X^{[t]}}} \tilde{x}^{(t+1)}\|. \tag{29}$$

Therefore, we first bound the norm of $\tilde{x}^{(t+1)}$ which turns out to be $\|\tilde{x}^{(t+1)}\| = \tilde{\Theta} \left(\frac{\sqrt{k}}{d} \right)$ as argued in the following.

Lower bound: The lower bound on $\|\tilde{x}^{(t+1)}\|$ simply follows from the term $u^{(t+1)}$, which is an independent random Gaussian.

Claim 4. *If $t \leq \frac{d}{10}$, then we have whp*

$$\|\tilde{x}^{(t+1)}\| \geq \frac{\delta'_t \sqrt{k}}{4d}.$$

Proof: We have

$$\|\tilde{x}^{(t+1)}\| \geq \|P_{\text{span}(X^{[t]}, U^{[t]}, a_1)^\perp} \tilde{x}^{(t+1)}\| = \|P_{\text{span}(X^{[t]}, U^{[t]}, a_1)^\perp} u^{(t+1)}\|.$$

Note that the equality is concluded from expansion of $\tilde{x}^{(t+1)}$ in (28) where all the components of $\tilde{x}^{(t+1)}$ in the subspace $\text{span}(X^{[t]}, U^{[t]}, a_1)^\perp$ is represented by $u^{(t+1)}$. The vector $P_{\text{span}(X^{[t]}, U^{[t]}, a_1)^\perp} u^{(t+1)}$ is the projection of a random Gaussian vector $u^{(t+1)}$ in to a subspace of dimation $d - o(d)$. Hence it is still a random Gaussian vector with expected square norm larger than $\frac{\delta'_t{}^2}{2} \frac{k}{d^2}$. By Lemma 16, with high probability the desired bound holds. \square

Upper bound: The upper bound is argued in the following claim.

Claim 5. *We have either*

$$\langle x^{(t+1)}, a_1 \rangle \geq 1 - \gamma,$$

for some constant $\gamma > 0$ or

$$\|\tilde{x}^{(t+1)}\| \leq \tilde{O}\left(\frac{\sqrt{k}}{d}\right).$$

Proof: Let $\tilde{x}^{(t+1)}$ in (28) be written as $\tilde{x}^{(t+1)} = z + (y_1^{(t)})^2 a_1$ where vector $z \in \mathbb{R}^d$ represents all the other terms in the expansion. The analysis is done under two cases 1) $(y_1^{(t)})^2 \geq \frac{2}{\gamma} \|z\|$ and 2) $(y_1^{(t)})^2 < \frac{2}{\gamma} \|z\|$ for some constant $\gamma > 0$. Note that the left hand side is the norm of $(y_1^{(t)})^2 a_1$ since $\|a_1\| = 1$, and in addition $(y_1^{(t)})^2 = \langle x^{(t)}, a_1 \rangle^2$.

Case 1 $\left((y_1^{(t)})^2 \geq \frac{2}{\gamma} \|z\|\right)$: For the $x^{(t+1)} := \tilde{x}^{(t+1)} / \|\tilde{x}^{(t+1)}\|$, we have

$$\begin{aligned} \langle x^{(t+1)}, a_1 \rangle &= \frac{1}{\|z + (y_1^{(t)})^2 a_1\|} \langle z + (y_1^{(t)})^2 a_1, a_1 \rangle \\ &\geq \frac{1}{\|z\| + (y_1^{(t)})^2} \left[(y_1^{(t)})^2 - \|z\| \right] \\ &\geq \frac{1 - \frac{\gamma}{2}}{1 + \frac{\gamma}{2}} \geq 1 - \gamma, \end{aligned}$$

where triangle and Cauchy-Schwartz inequality are used in the first bound, and the second inequality is concluded from assumption $(y_1^{(t)})^2 \geq \frac{2}{\gamma} \|z\|$.

Case 2 $\left((y_1^{(t)})^2 < \frac{2}{\gamma} \|z\|\right)$: We exploit the induction hypothesis to bound the norm of all the terms in the expansion of $\tilde{x}^{(t+1)}$ in (28).

For the terms involving $u^{(i+1)}$, $i \in [t]$, we have $\|u^{(i+1)}\| \leq 2\Delta'_i \frac{\sqrt{k}}{d}$ from Hypothesis 4 and the argument for $\|u^{(t+1)}\|$. In addition, for $u^{(i+1)}$, $i \in [t-1]$, the coefficient is bounded as

$$\frac{\langle P_{\perp_{W^{[i-1]}}} w^{(i)}, w^{(t)} \rangle}{\|P_{\perp_{W^{[i-1]}}} w^{(i)}\|^2} \leq \frac{\|w^{(t)}\|}{\|P_{\perp_{W^{[i-1]}}} w^{(i)}\|} \leq \frac{\Delta'_t}{\delta'_i}, \quad (30)$$

where Cauchy-Schwartz inequality is used in the first bound, and the bound in Hypothesis 2 and (27) are exploited in the last inequality. Therefore, the total contribution from terms involving $u^{(i+1)}$ in $\|\tilde{x}^{(t+1)}\|$ is bounded by $\frac{2(t-1)\Delta'_t{}^2 \sqrt{k}}{\delta'_t}$.

For the terms involving $P_{\perp_{X^{[i-1]}}} x^{(i)}$, $i \in [t]$, we have $\|P_{\perp_{X^{[i-1]}}} x^{(i)}\| \leq 1$, but the coefficient $\langle v^{(i)}, w^{(t)} \rangle$ needs further analysis to be bounded which is done in Lemma 12 saying $|\langle v^{(i)}, w^{(t)} \rangle| \leq \tilde{O}\left(\frac{\sqrt{k}}{d}\right)$. This implies that the total contribution from terms involving $P_{\perp_{X^{[i-1]}}} x^{(i)}$ in $\|\tilde{x}^{(t+1)}\|$ is bounded by $\tilde{O}\left(\frac{\sqrt{k}}{d}\right)$.

Combining the above bounds and considering the assumption that the norm of $(y_1^{(t)})^2 a_1$ in the expansion of $\tilde{x}^{(t+1)}$ is dominated by the norm of other terms argued above, the proof is complete concluding that $\|\tilde{x}^{(t+1)}\| \leq \tilde{O}\left(\frac{\sqrt{k}}{d}\right)$. \square

Lemma 12. *Under the induction hypothesis (up to update step $\tilde{x}^{(t+1)} := A(y^{(t)})^{*2}$ at iteration t), we have for $i \in [t]$,*

$$|\langle v^{(i)}, w^{(t)} \rangle| \leq O\left(t^3 \frac{(\Delta'_{t-1})^4}{(\delta'_{t-1})^4 \delta'_t} (\log d) \frac{\sqrt{k}}{d}\right) = \tilde{O}\left(\frac{\sqrt{k}}{d}\right).$$

Using (29) and the fact that $\|\tilde{x}^{(t+1)}\| = \tilde{\Theta}\left(\frac{\sqrt{k}}{d}\right)$, we have

$$\|P_{\perp_{X^{[t]}}} x^{(t+1)}\| \geq \tilde{\Theta}\left(\frac{d}{\sqrt{k}}\right) \|P_{\text{span}(X^{[t]}, U^{[t]}, a_1)^\perp} u^{(t+1)}\| \geq \frac{\delta'_t}{4},$$

where the bound $\|P_{\text{span}(X^{[t]}, U^{[t]}, a_1)^\perp} u^{(t+1)}\| \geq \frac{\delta'_t \sqrt{k}}{4}$ is also used. This finishes the proof that Hypothesis 1 holds.

Hypothesis 3

Finally we prove Hypothesis 3 at iteration t given earlier induction hypothesis. The first part of the hypothesis is proved in the following claim.

Claim 6. *We have*

$$|\langle a_1, x^{(t+1)} \rangle| \in [\delta_{t+1}^*, \Delta_{t+1}^*] d^{\beta 2^t} \frac{\sqrt{k}}{d}.$$

Proof: We first show the correlation bound on the unnormalized version as $\langle a_1, \tilde{x}^{(t+1)} \rangle$. Looking at the expansion of $\tilde{x}^{(t+1)}$ in (28), the correlation $\langle a_1, \tilde{x}^{(t+1)} \rangle$ involves three types of terms emerging from $(y_1^{(t)})^2 a_1$, $u^{(i+1)}$ and $P_{\perp_{X^{(i-1)}}} x^{(i)}$. In the following, we argue the correlation from each of these terms where we observe that the correlation is dominated by the term $(y_1^{(t)})^2 a_1$, and the rest of terms contribute much smaller amount.

For the term $(y_1^{(t)})^2 a_1$, we have

$$\langle a_1, (y_1^{(t)})^2 a_1 \rangle = (y_1^{(t)})^2 \in [(\delta_t^*)^2, (\Delta_t^*)^2] d^{\beta 2^t} \frac{k}{d^2},$$

where the last part exploits induction Hypothesis 3 in the previous iteration.

For the terms involving $u^{(i+1)}$, these vectors are random Gaussian vectors in a subspace (with dimension $\Omega(d)$), and therefore, we have with high probability

$$\langle a_1, u^{(i+1)} \rangle \leq \mathbb{E}[\|u^{(i+1)}\|] \cdot O\left(\frac{\log d}{\sqrt{d}}\right) \leq \tilde{O}\left(\frac{\sqrt{k}}{d\sqrt{d}}\right) \leq \tilde{O}\left(\frac{k}{d^2}\right),$$

where the correlation bound between two independent random Gaussian vectors in $\Omega(d)$ -dimension is used in the first inequality¹², Hypothesis 4 is exploited in the second inequality, and finally last inequality is from assumption $k > d$. In addition, the coefficient associated with $u^{(i+1)}$ is bounded by Δ'_t/δ'_i argued in (30). Hence, the total contribution from terms involving $u^{(i+1)}$ in $\langle \tilde{x}^{(t+1)}, a_1 \rangle$ is bounded by $\tilde{O}\left(\frac{k}{d^2}\right)$.

For the terms involving $P_{\perp_{X^{[i-1]}}} x^{(i)}$, by Hypothesis 3 we have

$$\langle a_1, P_{\perp_{X^{[i-1]}}} x^{(i)} \rangle \leq \Delta_i^* d^{\beta 2^{i-1}} \frac{\sqrt{k}}{d}.$$

In addition, the associated coefficient is bounded by $\tilde{O}\left(\frac{\sqrt{k}}{d}\right)$ from Lemma 12. Hence, the total contribution from terms involving $P_{\perp_{X^{[i-1]}}} x^{(i)}$ in $\langle \tilde{x}^{(t+1)}, a_1 \rangle$ is bounded by $\tilde{O}\left(d^{\beta 2^{t-1}} \frac{k}{d^2}\right)$.

Combining the above bounds implies

$$|\langle a_1, \tilde{x}^{(t+1)} \rangle| \leq \tilde{O}\left(d^{\beta 2^t} \frac{k}{d^2}\right).$$

Finally, using the bound on the norm of $\tilde{x}^{(t+1)}$ argued as $\|\tilde{x}^{(t+1)}\| = \tilde{\Theta}\left(\frac{\sqrt{k}}{d}\right)$ finishes the proof. \square

To prove the last part of Hypothesis 3, we use the following lemma which is very similar to Lemma 11.

Lemma 13. *Suppose R is a subspace in \mathbb{R}^d of dimension t' , such that there is a basis $\{r_1, \dots, r_{t'}\}$ with $|\langle r_i, a_1 \rangle| \leq \Delta$ and $\|r_i\| = 1$. Let $p \in \mathbb{R}^d$ be an arbitrary vector, then*

$$|\langle P_{\perp_R} p, a_1 \rangle| \leq |\langle p, a_1 \rangle| + \|p\| \Delta \sqrt{t'}.$$

We apply this lemma with $R = \text{span}(X^{[t]})$, and the basis is $P_{\perp_{X^{[i-1]}}} X^{(i)} / \|P_{\perp_{X^{[i-1]}}} X^{(i)}\|$. By induction hypothesis Δ in the lemma is at most $\Delta_t^* d^{\beta 2^t} \sqrt{k}/d$, let $v = x^{(t+1)}$ then this gives the desired bound.

Let $R = \text{span}(X^{[t]})$. Then the vectors $P_{\perp_{X^{[i-1]}}} x^{(i)} / \|P_{\perp_{X^{[i-1]}}} x^{(i)}\|$, $i \in [t]$ form a basis for subspace R , and we know from Hypotheses 1 and 3 that the correlation between these basis vectors and a_1 is bounded by $\Delta := \Delta_t^* d^{\beta 2^{t-1}} \frac{\sqrt{k}}{d}$. Applying above lemma, we have

$$|\langle P_{\perp_{X^{[t]}}} x^{(t+1)}, a_1 \rangle| \leq |\langle x^{(t+1)}, a_1 \rangle| + \Delta \sqrt{t} \leq \Delta_{t+1}^* d^{\beta 2^t} \frac{\sqrt{k}}{d},$$

where the last inequality uses the first part of Hypothesis 3 proved earlier in this section. Note that Δ_{t+1}^* is a new polylog factor here.

¹²For two independent random Gaussian vectors $p, q \in \mathbb{R}^d$, we have with high probability $\langle p, q \rangle \leq \mathbb{E}[\|p\|] \cdot \mathbb{E}[\|q\|] \cdot O\left(\frac{\log d}{\sqrt{d}}\right)$.

A.3 Growth rate of $\delta_t, \delta'_t, \Delta'_t, \delta_t^*, \Delta_t^*$

We know that if the number of iterations t is a constant, then the δ and Δ parameters (i.e., $\delta_t, \delta'_t, \Delta'_t, \delta_t^*, \Delta_t^*$) in the induction hypothesis are bounded by polylog factors of d . Here, we show that these parameters can be still bounded even when the number of steps is slightly larger than a constant. Let

$$R_t := \max\{1/\delta_t, \Delta'_{t-1}/\delta'_{t-1}, \Delta_t^*/\delta_t^*\}.$$

We know $R_1 = 1$, and by the inductive step analysis we have the following polynomial recursion property.

Claim 7. $R_{t+1} = \text{poly}(R_t, t, \log d)$.

This claim follows from the proof of inductive step, where in every step the δ and Δ parameters are bounded by polynomial functions of previous δ 's (Δ 's), t , and $\log d$.

We now solve this recursion as follows.

Lemma 14. *Suppose $R_{t+1} \leq c_0 R_t^{c_1} t^{c_2} (\log d)^{c_3}$ where c_0, c_1, c_2, c_3 are positive constants, and we know $R_1 = 1$. Then*

$$R_t \leq (\log d)^{2^{c_4 t}},$$

for some constant $c_4 > 0$ depending on c_0, c_1, c_2, c_3 .

Proof: Without loss of generality assume $c_0 \geq 1, c_2 \geq 1, c_3 \geq 1$, and $R_1 \geq \log d$. Given these assumptions, we have $R_t \geq \max\{c_0, t, \log d\}$, for $t \geq 1$. Applying this to the assumption $R_{t+1} \leq c_0 R_t^{c_1} t^{c_2} (\log d)^{c_3}$, we have

$$R_{t+1} \leq R_t^{1+c_1+c_2+c_3}. \quad (31)$$

Pick some $q > 0$ such that $R_1 \leq (\log d)^{2^q}$, and pick some

$$c_4 \geq \max\{q, \log_2(1 + c_1 + c_2 + c_3)\}.$$

Now we prove the result by the induction argument. Since $c_4 \geq q$, the basis of induction holds for R_1 . As the inductive step, suppose $R_t \leq (\log d)^{2^{c_4 t}}$. Applying this to (31), we have

$$R_{t+1} \leq (\log d)^{(1+c_1+c_2+c_3)2^{c_4 t}} \leq (\log d)^{2^{c_4(t+1)}},$$

where $2^{c_4} \geq (1 + c_1 + c_2 + c_3)$ is used in the last inequality. This finishes the inductive step and the result is proved. \square

Using the above bound, we show in the following corollary that the δ and Δ parameters in the induction hypothesis are bounded by polylog factors of d even if the number of steps t goes up to $c \log \log d$ for small enough constant c . In addition, we show that if $\beta \geq (\log d)^{-c_5}$ for some constant $c_5 > 0$, then the power method converges to a point $x^{(t)}$ which is constant close to the true component.

Corollary 2. *There exists a universal constant $c_5 > 0$ such that if*

$$\beta \geq (\log d)^{-c_5},$$

and the initial correlation is lower bounded by $d^{\beta} \frac{\sqrt{k}}{d}$ (see (13)), then with high probability the power method gets to a point that is constant close to the true component in $\Theta(\log \log d)$ number of steps.

Proof: Pick the number of steps to be $t = (\log \log d)/2c_4$, where c_4 is the constant in Lemma 14. Then, from Lemma 14 we have

$$R_t \leq (\log d)^{\sqrt{\log d}} \leq o(d),$$

where the last inequality can be shown by taking the log of both sides. This says that the analysis of inductive step still holds after such number of iterations.

Finally, by progress bound in (20), we can see that if $\beta \geq (\log d)^{-c_5}$, then the power method converges to a point $x^{(t)}$ which is constant close to the true component. \square

B Auxiliary Lemmas for Induction Argument

In this section we prove the lemmas used in arguing inductive step in Appendix A.2.

We first introduce the following lemma proposing a lower bound on the singular value of product of matrices.

Lemma 15 (Merikoski and Kumar 2004). *Let C and D be $k \times k$ matrices. If $1 \leq i \leq k$ and $1 \leq l \leq k - i + 1$, then*

$$\sigma_i(CD) \geq \sigma_{i+l-1}(C) \cdot \sigma_{k-l+1}(D),$$

where $\sigma_j(C)$ denotes the j -th singular value (in decreasing order) of matrix C .

B.1 Properties of random Gaussian vectors

We start with some basic properties of random Gaussian vectors. First as a simple fact, the norm of a random Gaussian vector is concentrated as follows which is proved via simple concentration inequalities.

Lemma 16. *Let $z \in \mathbb{R}^d$ be a random Gaussian vector with $\mathbb{E}[zz^\top] = \frac{1}{d}I$. Then we have with high probability $\frac{1}{2} \leq \|z\| \leq 2$.*

Next we show the ℓ_∞ norm of a Gaussian vector is small, even if it is projected on some subspace.

Lemma 17. *Let R be any linear subspace in \mathbb{R}^d and $z \in \mathbb{R}^d$ be a random Gaussian vector with $\mathbb{E}[zz^\top] = \frac{1}{d}I$. Then we have with high probability $\|P_{\perp R}z\|_\infty \leq \frac{\log d}{\sqrt{d}}$.*

Proof: Since $P_{\perp R}$ is a projection matrix, in particular the norm of its columns is bounded by 1. Hence, each entry of $P_{\perp R}z$ is a Gaussian random variable with variance bounded by $\frac{1}{d}$ implying that with high probability the absolute value of each coordinate is smaller than $\frac{\log d}{\sqrt{d}}$. Finally, the desired ℓ_∞ norm bound is argued by applying union bound. \square

We can also show that most of the entries are of size at least $\frac{1}{\sqrt{d}}$.

Lemma 18. *Let R be any linear subspace in \mathbb{R}^d with dimension $t \leq \frac{d}{16(\log d)^2}$ and $z \in \mathbb{R}^d$ be a random Gaussian vector with $\mathbb{E}[zz^\top] = \frac{1}{d}I$. Then we have with high probability at least 1/2 of the entries $i \in [d]$ satisfy $|(P_{\perp R}z)_i| \geq \frac{1}{4\sqrt{d}}$.*

Proof: Since the entries of z are independent Gaussian random variables with standard deviation $\frac{1}{\sqrt{d}}$, we know with high probability at least 1/2 of the entries have absolute value larger than $\frac{1}{2\sqrt{d}}$. On the other hand, $P_R z$ is also a random Gaussian vector with expected square norm bounded by

$$\mathbb{E}[\|P_R z\|^2] \leq \frac{\mathbb{E}[\|z\|^2]}{16(\log d)^2} = \frac{1}{16(\log d)^2},$$

where the assumption on the dimension of subspace R is used in the inequality. By Lemma 17 we know with high probability entries of $P_R z$ are bounded by $1/4\sqrt{d}$. Now $P_{\perp R} z = z - P_R z$ must have at least $1/2$ of the entries with absolute value larger than $1/4\sqrt{d}$. \square

Using the above lemmas, we can prove Lemma 10.

Lemma 10 (Restated). *Suppose R and R' are two subspaces in \mathbb{R}^k with dimension at most $t \leq \frac{k}{16(\log k)^2}$. Let $p \in \mathbb{R}^k$ be an arbitrary vector, $z \in \mathbb{R}^k$ be a uniformly random Gaussian vector in the space orthogonal to R , and finally $w = (p + z) * (p + z)$. Then with high probability, we have*

$$\|P_{\perp R'} w\| \geq \frac{\mathbb{E}[\|z\|^2]}{40\sqrt{k}}.$$

Proof: Let z, z' be two independent samples of z , and w, w' be the corresponding w vectors. We have

$$w - w' = (p + z) * (p + z) - (p + z') * (p + z') = (2p + z + z') * (z - z'). \quad (32)$$

By properties of Gaussian vectors, $z + z', z - z'$ are two *independent* random Gaussian vectors in the subspace orthogonal to R each with expected square norm $2\mathbb{E}[\|z\|^2]$. We use $z_1 := z + z'$ and $z_2 := z - z'$ to denote these two random Gaussian vectors.

Next, we show that with high probability

$$\|P_{\perp R'}(w - w')\| \geq \frac{\mathbb{E}[\|z\|^2]}{20\sqrt{k}}.$$

Note that this implies the result of lemma as follows. Suppose $\|P_{\perp R'} w\| < \frac{1}{40}\mathbb{E}[\|z\|^2]/\sqrt{k}$ with probability δ . Since w' is an independent sample, with probability δ^2 this bound holds for both w and w' . When this happens, we have $\|P_{\perp R'}(w - w')\| < \frac{1}{20}\mathbb{E}[\|z\|^2]/\sqrt{k}$ by triangle inequality. Since we showed δ^2 is negligible, δ is also negligible.

First we sample z_2 . Let $R'' = \text{span}(R', p * z_2)$. Then by expansion of $w - w'$ in (32), we have

$$\|P_{\perp R'}(w - w')\| = \|P_{\perp R'}(2(p * z_2) + (z_1 * z_2))\| \geq \|P_{\perp R''}(z_1 * z_2)\| = \|P_{\perp R''} \text{Diag}(z_2) P_{\perp R} z_1\|, \quad (33)$$

where the inequality is concluded by ignoring the component along $p * z_2$ direction. The last equality is from¹³ $u * v = \text{Diag}(u) \cdot v$ (for two vectors u and v), and the assumption that $z_1 = z + z'$ is in the subspace orthogonal to R . For the matrix $P_{\perp R''} \text{Diag}(z_2) P_{\perp R}$, we have¹⁴

$$\sigma_{k/4}(P_{\perp R''} \text{Diag}(z_2) P_{\perp R}) \geq \sigma_{k/2}(\text{Diag}(z_2)) \cdot \sigma_{7k/8}(P_{\perp R}) \cdot \sigma_{7k/8}(P_{\perp R''}) \geq \frac{\sqrt{\mathbb{E}[\|z\|^2]}}{4\sqrt{k}},$$

where the first inequality is from Lemma 15, and the last step is argued as follows. By Lemma 16, with high probability z_2 has square norm at least $\mathbb{E}[\|z_2\|^2]/2 = \mathbb{E}[\|z\|^2]$, and therefore, by Lemma 18 at least $k/2$ of its entries have absolute value larger than $\frac{1}{4}\sqrt{\mathbb{E}[\|z\|^2]}/\sqrt{k}$. Therefore, we can restrict attention to the space spanned by the $k/4$ top singular vectors. In addition, within this subspace we have with high probability $\|z_1\|^2 \geq \mathbb{E}[\|z\|^2]/8$, and hence,

$$\|P_{\perp R''} \text{Diag}(z_2) P_{\perp R} z_1\| \geq \frac{\mathbb{E}[\|z\|^2]}{20\sqrt{k}},$$

which finishes the proof by applying (33). \square

¹³For vector u , $\text{Diag}(u)$ denotes the diagonal matrix with u as its main diagonal.

¹⁴Recall that $\sigma_l(A)$ denotes the l -th singular value (in decreasing order) of matrix A .

B.2 Properties of projections

In this part we prove some basic properties of projections. Intuitively, if the whole subspace has small inner-product with some vector, then the projection of an arbitrary vector to the orthogonal subspace should not change the inner-product with that particular vector by too much. This is what we require in Lemma 13.

Lemma 13 (Restated). *Suppose R is a subspace in \mathbb{R}^d of dimension t' , such that there is a basis $\{r_1, \dots, r_{t'}\}$ with $|\langle r_i, a_1 \rangle| \leq \Delta$ and $\|r_i\| = 1$. Let $p \in \mathbb{R}^d$ be an arbitrary vector, then*

$$|\langle P_{\perp R} p, a_1 \rangle| \leq |\langle p, a_1 \rangle| + \|p\| \Delta \sqrt{t'}.$$

Proof: We have $P_{\perp R} p = p - \sum_{i=1}^{t'} \langle p, r_i \rangle r_i$, and therefore

$$\begin{aligned} |\langle P_{\perp R} p, a_1 \rangle| &\leq |\langle p, a_1 \rangle| + \sum_{i=1}^{t'} |\langle p, r_i \rangle \langle a_1, r_i \rangle| \\ &\leq |\langle p, a_1 \rangle| + \Delta \sum_{i=1}^{t'} |\langle p, r_i \rangle| \\ &\leq |\langle p, a_1 \rangle| + \Delta \sqrt{t' \sum_{i=1}^{t'} \langle p, r_i \rangle^2} \\ &\leq |\langle p, a_1 \rangle| + \Delta \|p\| \sqrt{t'}. \end{aligned}$$

The first step is triangle inequality and the third is Cauchy-Schwartz. □

Lemma 11 is very similar.

Lemma 11 (Restated). *Suppose R is a subspace in \mathbb{R}^k of dimension t' , such that there is a basis $\{r_1, \dots, r_{t'}\}$ with $\|r_i\|_{\infty} \leq \frac{\Delta}{\sqrt{k}}$ and $\|r_i\| = 1$. Let $p \in \mathbb{R}^k$ be an arbitrary vector, then*

$$\|P_{\perp R} p\|_{\infty} \leq \|p\|_{\infty} + \|p\| \Delta \frac{\sqrt{t'}}{\sqrt{k}}.$$

This lemma essentially follows from Lemma 13, because ℓ_{∞} norm is just the maximum inner-product to a basis vector. More specifically, the above lemma is applied for all $a_1 = e_j, j \in [k]$, where e_j denotes the j -th basis vector in \mathbb{R}^k .

B.3 Bounding correlation between v and w

We are only left with Lemma 12. The main difficulty in proving this lemma is that the later steps are dependent on the previous steps. In the proof we show the dependency is bounded and in fact we can treat them as independent.

Lemma 12 (Restated). *Under the induction hypothesis (up to update step $\tilde{x}^{(t+1)} := A(y^{(t)})^{*2}$ at iteration t), we have for $i \in [t]$,*

$$|\langle v^{(i)}, w^{(t)} \rangle| \leq O \left(t^3 \frac{(\Delta'_{t-1})^4}{(\delta'_{t-1})^4 \delta_t^2} (\log d) \frac{\sqrt{k}}{d} \right) = \tilde{O} \left(\frac{\sqrt{k}}{d} \right).$$

Proof: Recall $w^{(t)} = y_{-1}^{(t)} * y_{-1}^{(t)}$, and $y_{-1}^{(t)}$ is specified in (24). We now expand the Hadamard product in $w^{(t)}$ and bound all the resulting $O(t^2)$ terms.

The first type of terms has the form $\langle v^{(i)}, P_{\perp_{W^{[i_1-1]}}} w^{(i_1)} * P_{\perp_{W^{[i_2-1]}}} w^{(i_2)} \rangle$, which can be bounded as

$$\begin{aligned} \langle v^{(i)}, P_{\perp_{W^{[i_1-1]}}} w^{(i_1)} * P_{\perp_{W^{[i_2-1]}}} w^{(i_2)} \rangle &\leq k \cdot \|v^{(i)}\|_{\infty} \cdot \|P_{\perp_{W^{[i_1-1]}}} w^{(i_1)} * P_{\perp_{W^{[i_2-1]}}} w^{(i_2)}\|_{\infty} \\ &\leq 2k \frac{\log d (\Delta'_{t-1})^2}{\sqrt{d} d^2}, \end{aligned}$$

where $\|v^{(i)}\|_{\infty}$ is bounded by Lemma 17, and ℓ_{∞} norm of other vector is bounded by induction Hypothesis 2. In addition, the corresponding coefficient is bounded by (see (26), and note that both $i_1, i_2 < t$)

$$\frac{4(\Delta'_{t-1})^2 d^2}{(\delta'_{t-1})^4 k}.$$

Hence, the total contribution from such terms is bounded by

$$8t^2 \frac{(\Delta'_{t-1})^4 \log d}{(\delta'_{t-1})^4 \sqrt{d}}. \quad (34a)$$

The second type of terms has the form $\langle v^{(i)}, P_{\perp_{W^{[i_1-1]}}} w^{(i_1)} * v^{(i_2)} \rangle = \langle v^{(i)} * v^{(i_2)}, P_{\perp_{W^{[i_1-1]}}} w^{(i_1)} \rangle$, which can be bounded as

$$\|P_{\perp_{W^{[i_1-1]}}} w^{(i_1)}\|_{\infty} \cdot \|v^{(i)} * v^{(i_2)}\|_1 \leq \|P_{\perp_{W^{[i_1-1]}}} w^{(i_1)}\|_{\infty} \cdot \frac{\|v^{(i)}\|^2 + \|v^{(i_2)}\|^2}{2} \leq 4\Delta'_{t-1} \frac{k}{d^2},$$

where the last inequality is concluded from Hypotheses 2 and 4. In addition, the corresponding coefficient is bounded by (see (25) and (26), and note that both $i_1, i_2 < t$)

$$\frac{2\Delta'_{t-1} d}{(\delta'_{t-1})^2 \delta_{t-1} \sqrt{k}}.$$

Hence, the total contribution from such terms is bounded by

$$8t^2 \frac{(\Delta'_{t-1})^2 \sqrt{k}}{\delta_{t-1} (\delta'_{t-1})^2 d}. \quad (34b)$$

The third type of terms has the form $\langle v^{(i)}, v^{(i_1)} * v^{(i_2)} \rangle$, with coefficient bounded by $1/\delta_{t-1}^2$ (see (25)). For bounding these inner products, we need to use the fact that they are random Gaussian vectors, however the main difficulty is that they are correlated (if $i > j$, then the subspace that $v^{(i)}$ is in that depends on $v^{(j)}$). To resolve this difficulty, we treat $v^{(i)} \in \mathbb{R}^{k-1}$ as projection of $n^{(i)} \in \mathbb{R}^{k-1}$ into subspace orthogonal to $W^{[t-1]}$, where $n^{(i)}$'s are *independent* Gaussian vectors in the full $k-1$ dimensional space. Independent of the ordering of i, i_1, i_2 , we have with high probability

$$\langle n^{(i)}, n^{(i_1)} * n^{(i_2)} \rangle \leq O\left(\frac{\sqrt{k}}{d\sqrt{d}}\right),$$

since it is a sum of $k - 1$ independent mean-0 entries each with variance $\frac{1}{d^3}$. On the other hand, from Hypothesis 4, we have $\mathbb{E}[\|v^{(i)}\|^2] \leq 4\frac{k}{d}$, and since vector $n^{(i)} - v^{(i)}$ is in the subspace $W^{[t-1]}$ with dimension t , we have

$$\mathbb{E}[\|n^{(i)} - v^{(i)}\|^2] \leq O\left(\frac{t}{k}\right) \cdot \frac{4k}{d} = O\left(\frac{t}{d}\right),$$

and therefore, we have with high probability $\|n^{(i)} - v^{(i)}\| \leq O(\sqrt{t/d})$ for all $i \in [t-1]$. Using this, the difference between $\langle n^{(i)}, n^{(i_1)} * n^{(i_2)} \rangle$ and $\langle v^{(i)}, v^{(i_1)} * v^{(i_2)} \rangle$ can be bounded as

$$|\langle n^{(i)}, n^{(i_1)} * n^{(i_2)} \rangle - \langle v^{(i)}, v^{(i_1)} * v^{(i_2)} \rangle| \leq O\left((\log k)t \frac{\sqrt{k}}{d\sqrt{d}}\right),$$

where the right hand side is the bound on the dominant term in the expansion of difference as

$$\begin{aligned} |\langle n^{(i)}, (n^{(i_1)} - v^{(i_1)}) * (n^{(i_2)} - v^{(i_2)}) \rangle| &\leq \|n^{(i)}\| \cdot \|(n^{(i_1)} - v^{(i_1)}) * (n^{(i_2)} - v^{(i_2)})\| \\ &\leq O\left((\log k)\sqrt{\frac{k}{d}}\right) \cdot O\left(\frac{t}{d}\right) \\ &= O\left((\log k)t \frac{\sqrt{k}}{d\sqrt{d}}\right). \end{aligned}$$

Here, the first inequality is the Cauchy-Schwartz, and the second inequality is from bound on the norm of random Gaussian vector $n^{(i)}$, and the bound on the norm of difference vectors $n^{(i_1)} - v^{(i_1)}$ stated earlier. Hence, the total contribution from such terms is bounded by

$$O\left(t^3 \frac{\log k}{\delta_{t-1}^2} \frac{\sqrt{k}}{d\sqrt{d}}\right). \quad (34c)$$

Taking the sum of all the terms in (34a)-(34c) gives the desired bound. □

C Additional Arguments for Noise Analysis

Proof of Lemma 9: We prove this by an induction argument.

Basis of induction: For $t = 1$, $x^{(1)}$ is the initialization vector and thus, $\xi^{(1)} = 0$. Hence, the proposed bound holds for the basis of induction $t = 1$.

Inductive step: Assuming the inductive hypothesis holds for step t , we prove it also holds for step $t + 1$. We have

$$\begin{aligned} x^{(t+1)} + \xi^{(t+1)} &= \text{Norm}\left(\hat{T}(x^{(t)} + \xi^{(t)}, x^{(t)} + \xi^{(t)}, I)\right) \\ &= \text{Norm}\left(T(x^{(t)}, x^{(t)}, I) + 2T(x^{(t)}, \xi^{(t)}, I) + T(\xi^{(t)}, \xi^{(t)}, I) + E(\hat{x}^{(t)}, \hat{x}^{(t)}, I)\right). \quad (35) \end{aligned}$$

The first term $T(x^{(t)}, x^{(t)}, I)$ corresponds to the main signal; recall that $x^{(t+1)} = \text{Norm}(T(x^{(t)}, x^{(t)}, I))$ in the noiseless setting, where the unnormalized version $\tilde{x}^{(t+1)} := T(x^{(t)}, x^{(t)}, I)$ has norm at least $\tilde{\Omega}(\sqrt{k}/d)$ which is argued in the induction argument for Hypothesis 1. We now bound the desired property of noise terms in the above expansion.

For the second term, we break it into two terms as

$$2T(x^{(t)}, \xi^{(t)}, I) = 2\langle x^{(t)}, a_1 \rangle \langle \xi^{(t)}, a_1 \rangle a_1 + 2T'(x^{(t)}, \xi^{(t)}, I) =: p + q,$$

where $T' := \sum_{j>1} a_j \otimes a_j \otimes a_j$. Here $p := 2\langle x^{(t)}, a_1 \rangle \langle \xi^{(t)}, a_1 \rangle a_1$ corresponds to the multilinear form from first component of T , and $q := 2T'(x^{(t)}, \xi^{(t)}, I)$ corresponds to the multilinear form from the rest of components.

For q , we apply Lemma 19. Note that since $\|x^{(t)}\|_{B^*} \leq \tilde{O}(1/\sqrt{d})$, we get an extra $1/\sqrt{d}$ factor in the bound provided by Lemma 19, and therefore we have

$$\|q\|_2 \leq \tilde{O}(\epsilon d^{\beta 2^{t-1}} \sqrt{k}/d),$$

where we also used the induction hypothesis $\|\xi^{(t)}\| \leq \tilde{O}(\epsilon d^{\beta 2^{t-1}})$.

For p , we have

$$\|p\| = 2|\langle x^{(t)}, a_1 \rangle| \cdot |\langle \xi^{(t)}, a_1 \rangle| \leq \tilde{O}\left(\epsilon d^{\beta 2^t} \sqrt{k}/d\right),$$

where the inequality is from the signal and noise induction hypotheses; see Equation (20) for the signal induction hypothesis.

The third term $T(\xi^{(t)}, \xi^{(t)}, I)$ has ℓ_2 norm bounded as

$$\|T(\xi^{(t)}, \xi^{(t)}, I)\| \leq \|T\| \|\xi^{(t)}\|^2 \leq \tilde{O}(d^{\beta 2^t} \epsilon^2) \leq \tilde{O}(\epsilon d^{\beta 2^t} \sqrt{k}/d),$$

where the first inequality uses the sub-multiplicative property, and the second inequality exploits the bounded norm of random tensor T as $\|T\| \leq O(1)$, and the induction hypothesis in t -th step. The final inequality uses the assumption $\epsilon < o(\sqrt{k}/d)$ in the lemma.

The fourth term $E(\hat{x}^{(t)}, \hat{x}^{(t)}, I)$ has ℓ_2 norm bounded by

$$\|E(\hat{x}^{(t)}, \hat{x}^{(t)}, I)\| \leq \|E\| \|\hat{x}^{(t)}\|^2 \leq \epsilon \sqrt{k}/d,$$

where we use the sub-multiplicative property in the first inequality, and the assumption on the norm of error tensor E in the lemma, and the fact that $\|\hat{x}^{(t)}\| = 1$ are exploited in the second inequality.

Summarizing the above calculations on different terms of the update in (35), the signal plus noise vector before normalization is

$$T(x^{(t)}, x^{(t)}, I) + 2T(x^{(t)}, \xi^{(t)}, I) + T(\xi^{(t)}, \xi^{(t)}, I) + E(\hat{x}^{(t)}, \hat{x}^{(t)}, I) =: \alpha x^{(t+1)} + z,$$

where α is a coefficient which is lower bounded as $\alpha \geq \tilde{\Omega}(\sqrt{k}/d)$. The vector z also satisfies

$$\|z\| \leq \tilde{O}(\epsilon d^{\beta 2^t} \sqrt{k}/d), \tag{36}$$

which is derived by combining the bounds we argued on the second, third and fourth terms.

Note that until the very last step we always have $d^{\beta 2^t} \leq o(d/\sqrt{k})$ (otherwise we are constantly close to the true component, and we are done). In this case the norm of z is negligible compared

to α since $\|z\| \leq o(\alpha)$, and thus, the normalization factor is equal to $\|\alpha x^{(t+1)} + z\| = \alpha(1 \pm o(1))$. Therefore, after the normalization, we have the noise vector $\xi^{(t+1)} = \alpha' x^{(t+1)} + \beta z$, where $|\alpha'| \leq \|z\|/\alpha \leq o(1)$ and $|\beta| \leq 2/\alpha \leq \tilde{O}(d/\sqrt{k})$, hence we know $\|\xi^{(t+1)}\| \leq \tilde{O}(\epsilon d^{\beta 2^t})$.

For the last step of the induction, the norm of $T(x^{(t)}, x^{(t)}, I)$ is also larger (it has norm $d^{\beta 2^t} k/d^2$, which is larger than \sqrt{k}/d for the last step). Since $\epsilon < o(\sqrt{k}/d)$ we still know the noise is negligible. \square

Lemma on the property of $\|\cdot\|_*$ norm defined in Definition 1:

Lemma 19. *Consider a random tensor $T = \sum_{j \in [k]} a_j \otimes a_j \otimes a_j$ where a_j 's are zero-mean random Gaussian with expected unit norm. Let $A \in \mathbb{R}^{d \times k}$ be the matrix $[a_1, \dots, a_k]$, $T' = \sum_{j>1} a_j \otimes a_j \otimes a_j$ and $B \in \mathbb{R}^{d \times (k-1)}$ be the matrix $[a_2, a_3, \dots, a_k]$. Then for any vectors u, v such that $\|u\|_{B^*} \leq 1$ and $\|v\|_2 \leq 1$, with high probability we have*

$$\|T'(u, v, I)\|_2 \leq \tilde{O}\left(\sqrt{k/d}\right).$$

Proof: We prove this lemma along similar ideas provided in the proof of Anandkumar et al. (2014b, Claim 1). Let η_j 's be independent random ± 1 variables with $\Pr[\eta_j = 1] = 1/2$. We rewrite tensor T' as

$$T' = \sum_{j>1} \eta_j a_j \otimes a_j \otimes a_j. \tag{37}$$

Since a_j 's are zero-mean random Gaussian vectors, we have $\eta_j a_j \sim a_j$, and thus, the new T' has the same distribution as the original one. We now first sample vectors a_j 's, and this already makes the norm $\|\cdot\|_{B^*}$ well-defined. In addition, the value of η_j 's does not change the singular values of A or B . Also note that since a_j 's are zero-mean random Gaussian vectors with expected norm 1, they also satisfy with high probability the incoherence condition such that $|\langle a_i, a_j \rangle| \leq \tilde{O}(1/\sqrt{d})$ for all $i \neq j$. Thus, we condition on all these fixed events, and the only remaining random variables are just the η_j 's.

The proposed statement in the lemma is equivalent to bounding

$$\sup_{\|u\|_{B^*}=1, \|v\|=\|w\|=1} |T'(u, v, w)|.$$

In order to bound it, we provide an ϵ -net argument. We construct an ϵ -net such that for any vector $u \in \mathbb{R}^d$ with unit $\|\cdot\|_{B^*}$ norm, there is a vector u' in the net such that $\|B^\top(u - u')\| \leq 1/k^2$. We also construct standard ϵ -net for vectors $u, w \in \mathbb{R}^d$ with unit ℓ_2 norm. By standard construction, this ϵ -net has size $\exp(O(d \log d))$. We now show that for all u in ϵ -net with unit $\|\cdot\|_{B^*}$ norm, and all v, w in ϵ -net with unit ℓ_2 norm, the desired bound $|T'(u, v, w)| \leq \tilde{O}\left(\sqrt{k/d}\right)$ holds with high probability. Then for the other vectors (u, v, w) not in the ϵ -net, the result follows from their closest points in the net.

Now for a fixed triple (u, v, w) in the ϵ -net, we have

$$T'(u, v, w) = \sum_{j>1} \eta_j \langle u, a_j \rangle \langle v, a_j \rangle \langle w, a_j \rangle,$$

which is a sum of independent random variables; recall that the randomness is from η_j 's, and a_j 's are already sampled and thus they are fixed here. We partition the above sum into *large* and *small* terms as $T'(u, v, w) = S_L + S_{L^c}$ such that the summation S_L is the sum of large terms including terms in set

$$L := \left\{ j \in \{2, 3, \dots, k\} : |\langle v, a_j \rangle| \geq \log d / \sqrt{d} \vee |\langle w, a_j \rangle| \geq \log d / \sqrt{d} \right\},$$

and the rest are the small terms forming S_{L^c} . Note that $|\langle u, a_j \rangle| \leq 1$ since $\|u\|_{B^*} = 1$.

Bounding $|S_{L^c}|$: Since the variables are bounded in this summation corresponding to small terms, we use Bernstein's inequality, and thus with probability at least $1 - \delta$, we have $|S_{L^c}| \leq \frac{\sqrt{k \log 1/\delta} \cdot \text{polylog } d}{d}$ for the fixed point in the ε -net. By choosing small enough $\delta = \exp(-Cd \log d)$ (where C is large enough constant), we can apply the union bound on the ε -net, and conclude that for all the vectors in the net, $|S_{L^c}|$ is smaller than $\tilde{O}(\sqrt{k/d})$ with high probability.

Bounding $|S_L|$: Since the columns of matrix B are random Gaussian vectors, it satisfies the RIP property with high probability (see Remark 3 in Anandkumar et al. (2014b) for the precise definition of RIP), and thus by the definition of RIP and Lemma 3 in Anandkumar et al. (2014b), we have $\|B_L\| \leq 2$ where B_L is the sub-columns of matrix B specified by set L .

We now have

$$|S_L| \leq \sum_{j \in L} |\langle u, a_j \rangle| \cdot |\langle v, a_j \rangle| \cdot |\langle w, a_j \rangle| \leq \sum_{j \in L} |\langle v, a_j \rangle| \cdot |\langle w, a_j \rangle| \leq \|B_L^\top v\| \cdot \|B_L^\top w\| \leq 4,$$

where the second step uses the fact that $|\langle u, a_j \rangle| \leq 1$, the third step exploits Cauchy-Schwartz inequality, and the last step uses bound $\|B_L\| \leq 2$. Notice that matrix B is already sampled before we do the ε -net argument, and therefore, we do not need to do union bound over all u, v, w for this event.

Since we assume the overcomplete regime $k \geq d$, the bound on $|S_{L^c}|$ is dominant which finishes the proof. □

References

- A. Anandkumar, D. Hsu, and S. M. Kakade. A Method of Moments for Mixture Models and Hidden Markov Models. In *Proc. of Conf. on Learning Theory*, June 2012.
- A. Anandkumar, R. Ge, and M. Janzamin. Learning Overcomplete Latent Variable Models through Tensor Methods. In *Proceedings of the Conference on Learning Theory (COLT)*, Paris, France, July 2015.
- Anima Anandkumar, Rong Ge, and Majid Janzamin. Guaranteed Non-Orthogonal Tensor Decomposition via Alternating Rank-1 Updates. *arXiv preprint arXiv:1402.5180*, Feb. 2014a.
- Anima Anandkumar, Rong Ge, and Majid Janzamin. Sample Complexity Analysis for Learning Overcomplete Latent Variable Models through Tensor Methods. *arXiv preprint arXiv:1408.0553*, Aug. 2014b.

- Animashree Anandkumar, Rong Ge, Daniel Hsu, Sham M. Kakade, and Matus Telgarsky. Tensor Decompositions for Learning Latent Variable Models. *Journal of Machine Learning Research*, 15:2773–2832, 2014c.
- Sanjeev Arora and Ravi Kannan. Learning mixtures of separated nonspherical gaussians. *The Annals of Applied Probability*, 15(1A):69–92, 2005.
- Mohsen Bayati and Andrea Montanari. The dynamics of message passing on dense graphs, with applications to compressed sensing. *arXiv preprint arXiv:1001.3448*, Jan. 2010.
- A. Bhaskara, M. Charikar, A. Moitra, and A. Vijayaraghavan. Smoothed analysis of tensor decompositions. *arXiv preprint arXiv:1311.3651*, 2013.
- Emmanuel J Candes and Terence Tao. Near-optimal signal recovery from random projections: Universal encoding strategies? *Information Theory, IEEE Transactions on*, 52(12):5406–5425, 2006.
- Dustin Cartwright and Bernd Sturmfels. The number of eigenvalues of a tensor. *Linear Algebra and its Applications*, 438(2):942–952, January 2013.
- Adam Coates, Andrew Y Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *International Conference on Artificial Intelligence and Statistics*, pages 215–223, 2011.
- Sanjoy Dasgupta. Learning mixtures of gaussians. In *FOCS*, 1999.
- L. De Lathauwer, J. Castaing, and J.-F. Cardoso. Fourth-order cumulant-based blind identification of underdetermined mixtures. *Signal Processing, IEEE Transactions on*, 55(6):2965–2973, 2007.
- Rong Ge and Tengyu Ma. Decomposing overcomplete 3rd order tensors using sum-of-squares algorithms. *arXiv preprint arXiv:1504.05287*, April 2015.
- N. Goyal, S. Vempala, and Y. Xiao. Fourier pca. *arXiv preprint arXiv:1306.5825*, 2013.
- Christopher J. Hillar and Lek-Heng Lim. Most tensor problems are NP hard. *Journal of ACM*, 60(6), November 2013.
- Roger A. Horn and Charles R. Johnson. *Matrix analysis*. Cambridge university press, 2012.
- D. Hsu and S. M. Kakade. Learning Mixtures of Spherical Gaussians: Moment Methods and Spectral Decompositions. *arXiv preprint arXiv:1206.5766*, 2012.
- Daniel Hsu and Sham M Kakade. Learning mixtures of spherical gaussians: moment methods and spectral decompositions. In *Proceedings of the 4th conference on Innovations in Theoretical Computer Science*, pages 11–20. ACM, 2013.
- A. T. Kalai, A. Moitra, and G. Valiant. Efficiently learning mixtures of two gaussians. In *STOC*, 2010.
- M. S. Lewicki and T. J. Sejnowski. Learning overcomplete representations. *Neural computation*, 12(2):337–365, 2000.

- Jorma K. Merikoski and Ravinder Kumar. Inequalities for spreads of matrix sums and products. *Applied Mathematics E-Notes*, 4:150–159, 2004.
- A. Moitra and G. Valiant. Settling the polynomial learnability of mixtures of gaussians. In *FOCS*, 2010.
- Karl Pearson. Contributions to the mathematical theory of evolution. *Philosophical Transactions of the Royal Society of London. A*, 186:343–414, 1895.
- S. Vempala and G. Wang. A spectral algorithm for learning mixtures of distributions. In *FOCS*, 2002.
- T. Zhang and G. Golub. Rank-one approximation to high order tensors. *SIAM Journal on Matrix Analysis and Applications*, 23:534–550, 2001.