

Power Iteration for Tensor PCA

Jiaoyang Huang

New York University, New York, NY

JH4427@NYU.EDU

Daniel Z. Huang

California Institute of Technology, Pasadena, CA

DZHUANG@CALTECH.EDU

Qing Yang

University of Science and Technology of China, China

YANGQ@USTC.EDU.CN

Guang Cheng

University of California, Los Angeles & Purdue University, Los Angeles, CA

GUANGCHENG@UCLA.EDU

Editor: Animashree Anandkumar

Abstract

In this paper, we study the power iteration algorithm for the asymmetric spiked tensor model, as introduced in Richard and Montanari (2014). We give necessary and sufficient conditions for the convergence of the power iteration algorithm. When the power iteration algorithm converges, for the rank one spiked tensor model, we show the estimators for the spike strength and linear functionals of the signal are asymptotically Gaussian; for the multi-rank spiked tensor model, we show the estimators are asymptotically mixtures of Gaussian. This new phenomenon is different from the spiked matrix model. Using these asymptotic results of our estimators, we construct valid and efficient confidence intervals for spike strengths and linear functionals of the signals.

Keywords: Spiked tensor, tensor PCA, power iteration, statistical inference

1. Introduction

High order arrays, or tensors have been actively considered in neuroimaging analysis, topic modeling, signal processing and recommendation system Frolov and Oseledets (2017); Comon (2014); Hackbusch (2012); Karatzoglou et al. (2010); Rendle and Schmidt-Thieme (2010); Zhou et al. (2013); Simony et al. (2016); Cichocki et al. (2015); Sidiropoulos et al. (2017). Setting the stage, imagine that the signal is in the form of a large symmetric low-rank k -th order tensor

$$\mathbf{X}^* = \sum_{j=1}^r \beta_j \mathbf{v}_j^{\otimes k} \in \otimes^k \mathbb{R}^n, \quad (1)$$

where r ($r \ll n$) represents the rank and β_j are the strength of the signals. Such low-rank tensor components appear in various applications, e.g. community detection Anandkumar et al. (2014a); Jing et al. (2020), moments estimation for latent variable models Hsu and Kakade (2013); Anandkumar et al. (2014b) and hypergraph matching Duchenne et al. (2011). Suppose that we do not have access to perfect measurements about the entries of this signal tensor. The observations $\mathbf{X} = \mathbf{X}^* + \mathbf{Z}$ are contaminated by a substantial amount of random noise (reflected by the random tensor \mathbf{Z} which has i.i.d. Gaussian entries

with mean 0 and variance $1/n$.) The aim is to perform reliable estimation and inference on the unseen signal tensor \mathbf{X}^* . In literature, this is the spiked tensor model, introduced in Richard and Montanari (2014).

In the special case, when $k = 2$, the above model reduces to the well-known “spiked matrix model” Johnstone (2001). In this setting it is known that there is an order 1 critical signal-to-noise ratio β_c , such that below β_c , it is information-theoretical impossible to detect the spikes, and above β_c , it is possible to detect the spikes by Principal Component Analysis (PCA). A body of work has quantified the behavior of PCA in this setting Johnstone (2001); Baik et al. (2005); Baik and Silverstein (2006); Paul (2007); Benaych-Georges and Nadakuditi (2012); Bai and Yao (2012); Johnstone and Lu (2009); Birnbaum et al. (2013); Cai et al. (2013); Ma et al. (2013); Vu et al. (2013); Cai et al. (2015); El Karoui et al. (2008); Ledoit et al. (2012); Donoho et al. (2018). We refer readers to the review articles Johnstone and Paul (2018) for more discussion and references to this and related lines of work

Tensor problems are far more than an extension of matrices. Not only the more involved structures and high-dimensionality, many concepts are not well defined Kolda and Bader (2009), e.g. eigenvalues and eigenvectors, and most tensor problems are NP-hard Hillar and Lim (2013). Despite a large body of work tackling the spiked tensor model, there are several fundamental yet unaddressed challenges that deserve further attention.

Computational Hardness. The same as the spiked matrix model, for spiked tensor model, there is an order 1 critical signal-to-noise ratio β_k (depending on the order k). Below β_k , it is information-theoretical impossible to detect the spikes, and above β_k , the maximum likelihood estimator is a distinguishing statistics Chen et al. (2019, 2018a); Lesieur et al. (2017); Perry et al. (2020); Jagannath et al. (2018). In the matrix setting the maximum likelihood estimator is the top eigenvector, which can be computed in polynomial time by, e.g., power iteration. However, for order $k \geq 3$ tensor, computing the maximum likelihood estimator is NP-hard in generic setting. In this setting, it is widely believed that there is a regime of signal-to-noise ratios for which it is information theoretically possible to recover the signal but there is no known algorithm to efficiently approximate it. In the pioneering work Richard and Montanari (2014), the algorithmic aspects of this model have been studied under the special setting when the rank $r = 1$ and the signal-to-noise-ratio is β . They showed that tensor power iteration with random initialization recovers the signal provided $\beta \gtrsim n^{(k-1)/2}$, and tensor unfolding recovers the signal provided $\beta \gtrsim n^{(\lceil k/2 \rceil - 1)/2}$. Based on heuristic arguments, they predicted that the necessary and sufficient condition for power iteration to succeed is $\beta \gtrsim n^{(k-2)/2}$, and for tensor unfolding is $\beta \gtrsim n^{(k-2)/4}$. Langevin dynamics and gradient descent were studied in Arous et al. (2020), and shown to recover the signal provided $\beta \gtrsim n^{(k-2)/2}$. Later the sharp threshold $\beta \gtrsim n^{(k-2)/4}$ was achieved using Sum-of-Squares algorithms Hopkins et al. (2015, 2016); Kim et al. (2017) and sophisticated iteration algorithms Luo et al. (2020); Zhang and Xia (2018); Han et al. (2020). The necessary part of this threshold still remains open, and its relation with hypergraphic planted clique problem was discussed in Luo and Zhang (2020a,b).

Statistical inferences. In many applications, it is often the case that the ultimate goal is not to characterize the L_2 or “bulk” behavior (e.g. the mean squared estimation error) of the signals, but rather to reason about the signals along a few preconceived yet important directions. In the example of community detecting for hypergraphs, the entries of the vector \mathbf{v} can represent different community memberships. The testing of whether

any two nodes belong to the same community is reduced to the hypothesis testing problem of whether the corresponding entries of \mathbf{v} are equal. These problems can be formulated as estimation and inference for linear functionals of a signal, namely, quantities of the form $\langle \mathbf{a}, \mathbf{v}_j \rangle$, $1 \leq j \leq r$ with a prescribed vector \mathbf{a} . A natural starting point is to plug in an estimator $\widehat{\mathbf{v}}_j$ of \mathbf{v}_j , i.e. the estimator $\langle \mathbf{a}, \widehat{\mathbf{v}}_j \rangle$. However, most prior works Richard and Montanari (2014); Hopkins et al. (2015, 2016); Kim et al. (2017); Luo et al. (2020); Zhang and Xia (2018) on spiked tensor models focuses on the L_2 risk analysis, which is often too coarse to give tight uncertainty bound for the plug-in estimator. To further complicate matters, there is often a bias issue surrounding the plug-in estimator. Addressing these issues calls for refined risk analysis of the algorithms.

1.1 Our Contributions

We consider the power iteration algorithm given by the following recursion

$$\mathbf{u}_0 = \mathbf{u}, \quad \mathbf{u}_{t+1} = \frac{\mathbf{X}[\mathbf{u}_t^{\otimes(k-1)}]}{\|\mathbf{X}[\mathbf{u}_t^{\otimes(k-1)}]\|_2} \quad (2)$$

where $\mathbf{u} \in \mathbb{R}^n$ with $\|\mathbf{u}\|_2 = 1$ is the initial vector, and $\mathbf{X}[\mathbf{v}^{\otimes(k-1)}] \in \mathbb{R}^n$ is the vector with i -th entry given by $\langle \mathbf{X}, \mathbf{e}_i \otimes \mathbf{v}^{\otimes(k-1)} \rangle$. The estimators are given by

$$\widehat{\mathbf{v}} = \mathbf{u}_T, \quad \widehat{\beta} = \langle \mathbf{X}, \widehat{\mathbf{v}}^{\otimes k} \rangle. \quad (3)$$

for some large T . In the worst case scenario, i.e. with random initialization, power iteration algorithm underperforms tensor unfolding. However, if extra information about the signals \mathbf{v}_j is available, power iteration algorithm with a warm start can be used to obtain a much better estimator. In fact this approach is commonly used to obtain refined estimators. In this paper, we study the convergence and statistical inference aspects of the power iteration algorithm. The main contributions of this paper are summarized below,

Convergence criterion. We give necessary and sufficient conditions for the convergence of the power iteration algorithm in finite time. In the rank one case $r = 1$, we show that the power iteration algorithm converges to the true signal \mathbf{v} in finite time, provided $|\beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-2}| \gg 1$ where \mathbf{u} is the initialization vector. In the complementary setting, if $|\beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-2}| \ll 1$, the output of the power iteration algorithm behaves like random Gaussian vectors, and has no correlations with the signal. With random initialization, i.e. \mathbf{u} is a uniformly random vector on the unit sphere, our results assert that the power iteration algorithm converges in finite time, if and only if $\beta \gtrsim n^{(k-2)/2}$, which verifies the prediction in Richard and Montanari (2014). This is analogous to the PCA of spiked matrix model, where power iteration recovers the top eigenvalue. However, the multi-rank spiked tensor model, i.e. $r \geq 2$, is different from multi-rank spiked matrix. The power iteration algorithm for multi-rank spiked tensor model is more sensitive to the initialization, i.e. the power iteration algorithm converges if $\max_j |\beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle^{k-2}| \gg 1$. In this case, it converges to \mathbf{v}_{j_*} with $j_* = \operatorname{argmax}_j |\beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle^{k-2}|$.

Statistical inference We consider the statistical inference problem for the spiked tensor model. We develop the limiting distributions of the above power iteration estimators. In the rank one case, above the threshold $|\beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-2}| \gg 1$, we show that our estimator $\langle \mathbf{a}, \widehat{\mathbf{v}} \rangle$

(modulo some global sign) admits the following first order approximation

$$\langle \mathbf{a}, \widehat{\mathbf{v}} \rangle \approx \left(1 - \frac{1}{2\beta^2}\right) \langle \mathbf{a}, \mathbf{v} \rangle + \frac{\langle \mathbf{a}^\perp, \boldsymbol{\xi} \rangle}{\beta},$$

where $\mathbf{a}^\perp = \mathbf{a} - \langle \mathbf{a}, \mathbf{v} \rangle \mathbf{v}$, and $\boldsymbol{\xi} = \mathbf{Z}[\mathbf{v}^{\otimes(k-1)}]$, is an n -dim vector, with each entry i.i.d. $\mathcal{N}(0, 1/n)$ Gaussian random variables. For multi-rank spiked tensor model, the output of power iteration algorithm depends on the angle between the initialization \mathbf{u} and the signals \mathbf{v}_j . We consider the case that the initialization \mathbf{u} is a uniformly random vector on the unit sphere. For such initialization, very interestingly, our estimator $\langle \mathbf{a}, \widehat{\mathbf{v}} \rangle$ is asymptotically a mixture of Gaussian, with modes at $\langle \mathbf{a}, \mathbf{v}_j \rangle$ and mixture weights depending on the signal strength β_j . Using these asymptotic results of our estimators, we construct valid and efficient confidence intervals for the linear functionals $\langle \mathbf{a}, \mathbf{v}_j \rangle$. In a concurrent work, Xia et al. (2020) studied the statistical inference for several low rank tensor models, and asymptotic distributions were established under some essential conditions on the signal-to-noise ratio.

1.2 Notations:

For a vector $\mathbf{v} \in \mathbb{R}^n$, we denote its i -th coordinate as $\mathbf{v}(i)$. We equate k -th order tensors in $\otimes^k \mathbb{R}^n$ with vectors of dimension n^k , i.e. $\boldsymbol{\tau} = (\tau_{i_1 i_2 \dots i_k})_{1 \leq i_1, i_2, \dots, i_k \leq n}$. For any two k -th order tensors $\boldsymbol{\tau}, \boldsymbol{\eta} \in \otimes^k \mathbb{R}^n$, we denote their inner product as $\langle \boldsymbol{\tau}, \boldsymbol{\eta} \rangle := \sum_{1 \leq i_1, i_2, \dots, i_k \leq n} \tau_{i_1 i_2 \dots i_k} \eta_{i_1 i_2 \dots i_k}$. A k -th order tensor can act on a $(k-1)$ -th order tensor, and return a vector: $\boldsymbol{\tau} \in \otimes^k \mathbb{R}^n$ and $\boldsymbol{\eta} \in \otimes^{k-1} \mathbb{R}^n$

$$\boldsymbol{\tau}[\boldsymbol{\eta}] \in \mathbb{R}^n, \quad \boldsymbol{\tau}[\boldsymbol{\eta}](i) = \langle \boldsymbol{\tau}, \mathbf{e}_i \otimes \boldsymbol{\eta} \rangle = \sum_{1 \leq i_1, \dots, i_{k-1} \leq n} \tau_{i i_1 i_2 \dots i_{k-1}} \eta_{i_1 i_2 \dots i_{k-1}}.$$

We denote the L_2 norm of a vector \mathbf{v} as $\|\mathbf{v}\|$. We use $\stackrel{d}{=}$ for the equality in law, and $\stackrel{d}{\rightarrow}$ for the convergence in law. We denote the index sets $\llbracket a, b \rrbracket = \{a, a+1, a+2, \dots, b\}$ and $\llbracket n \rrbracket = \{1, 2, 3, \dots, n\}$. We use C to represent large universal constant, and c a small universal constant, which may be different from line by line. We write that $X = O(Y)$ or $X \lesssim Y$, if there exists some universal constant such that $|X| \leq CY$. We write $X = o(Y)$, or $X \ll Y$ if the ratio $|X|/Y \rightarrow 0$ as n goes to infinity. We write $X \gtrsim Y$, if there exists a universal constant such that $X \geq C|Y|$. We write $X \asymp Y$ if there exist universal constants such that $cY \leq |X| \leq CY$. We say an event holds with high probability, if for there exists $c > 0$, and n large enough, the event holds with probability at least $1 - n^{-c \log n}$.

An outline of the paper is given as follows. In Section 2.1, we state our main results for the rank-one spiked tensor model. In particular, with general initialization a distributional result for the power iteration algorithm is developed. Section 2.2 investigates the general rank- r spiked tensor model. A similar distributional result is established with general initialization as in Section 2.1. While with uniformly distributed initialization over the unit sphere, we obtain a multinomial distribution which yields a mixture Gaussian. Numerical simulations are presented in Section 3. All proofs and technical details are deferred to the appendix.

2. Main Results

2.1 Rank one spiked tensor model

In this section, we state our main results for the rank-one spiked tensor model (corresponding to $r = 1$ in (1)):

$$\mathbf{X} = \beta \mathbf{v}^{\otimes k} + \mathbf{Z},$$

where

- $\mathbf{X} \in \otimes^k \mathbb{R}^n$ is the k -th order tensor observation.
- $\mathbf{Z} \in \otimes^k \mathbb{R}^n$ is a noise tensor. The entries of \mathbf{Z} are i.i.d. standard $\mathcal{N}(0, 1/n)$ Gaussian random variables.
- $\beta \in \mathbb{R}$ is the signal size.
- $\mathbf{v} \in \mathbb{R}^n$ is an unknown unit vector to be recovered.

We obtain a distributional result for the power iteration algorithm (2) with general initialization \mathbf{u} : when $|\beta|$ is above certain threshold, \mathbf{u}_t converges to \mathbf{v} , and the error is asymptotically Gaussian; when $|\beta|$ is below the same threshold, the algorithm does not converge.

Theorem 1 *Fix the initialization $\mathbf{u} \in \mathbb{R}^n$ with $\|\mathbf{u}\|_2 = 1$ and $\langle \mathbf{u}, \mathbf{v} \rangle \gtrsim 1/\sqrt{n}$. If $|\beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-2}| \geq n^\varepsilon$ with arbitrarily small $\varepsilon > 0$, then for any fixed unit vector $\mathbf{a} \in \mathbb{R}^n$, and*

$$T \geq 1 + \frac{1}{\varepsilon} \left(\frac{1}{2} + \frac{2 \log |\beta|}{\log n} \right),$$

with probability $1 - O(n^{-c(\log n)^2})$, the power iteration estimator $(\hat{\beta}, \hat{\mathbf{v}}) = (\mathbf{X}[\mathbf{u}_T^{\otimes k}], \mathbf{u}_T)$ from (3) satisfies

1. If k is odd, and $\beta > 0$, then $(\hat{\beta}, \hat{\mathbf{v}})$ recovers (β, \mathbf{v}) in the following sense,

$$\begin{aligned} \langle \mathbf{a}, \hat{\mathbf{v}} \rangle = \langle \mathbf{a}, \mathbf{u}_T \rangle &= \left(1 - \frac{1}{2\beta^2} \right) \langle \mathbf{a}, \mathbf{v} \rangle + \frac{\langle \mathbf{a}, \boldsymbol{\xi} \rangle - \langle \mathbf{a}, \mathbf{v} \rangle \langle \mathbf{v}, \boldsymbol{\xi} \rangle}{\beta} \\ &+ O \left(\frac{\log n}{\beta^2 \sqrt{n}} + \frac{(\log n)^{3/2}}{\beta^{3/2} n^{3/4}} + \frac{|\langle \mathbf{a}, \mathbf{v} \rangle|}{\beta^4} \right), \end{aligned} \quad (4)$$

where $\boldsymbol{\xi} = \mathbf{Z}[\mathbf{v}^{\otimes(k-1)}]$, is an n -dim vector, with each entry i.i.d. $\mathcal{N}(0, 1/n)$ Gaussian random variable. And

$$\hat{\beta} = \mathbf{X}[\mathbf{u}_T^{\otimes k}] = \beta + \langle \boldsymbol{\xi}, \mathbf{v} \rangle - \frac{k/2 - 1}{\beta} + O \left(\frac{\log n}{|\beta| \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta|^{1/2} n^{3/4}} + \frac{1}{|\beta|^3} \right). \quad (5)$$

2. If k is odd, and $\beta < 0$ then $(\hat{\beta}, \hat{\mathbf{v}})$ recovers $(-\beta, -\mathbf{v})$, in the sense of (4) and (5) with (β, \mathbf{v}) replaced by $(-\beta, -\mathbf{v})$

3. If k is even, and $\beta > 0$, then $(\widehat{\beta}, \widehat{\mathbf{v}})$ recovers $(\beta, \text{sgn}(\langle \mathbf{u}, \mathbf{v} \rangle) \mathbf{v})$, in the sense of (4) and (5) with (β, \mathbf{v}) replaced by $(\beta, \text{sgn}(\langle \mathbf{u}, \mathbf{v} \rangle) \mathbf{v})$;
4. If k is even, and $\beta < 0$, then $(\widehat{\beta}, \widehat{\mathbf{v}})$ recovers $(\beta, (-1)^T \text{sgn}(\langle \mathbf{u}, \mathbf{v} \rangle) \mathbf{v})$, in the sense of (4) and (5) with (β, \mathbf{v}) replaced by $(\beta, (-1)^T \text{sgn}(\langle \mathbf{u}, \mathbf{v} \rangle) \mathbf{v})$.

Theorem 2 Fix the initialization $\mathbf{u} \in \mathbb{R}^n$ with $\|\mathbf{u}\|_2 = 1$. If $|\beta| \geq n^\varepsilon$ and $|\beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-2}| \leq n^{-\varepsilon}$ with arbitrarily small $\varepsilon > 0$, then \mathbf{u}_t does not converge to $\pm \mathbf{v}$, and \mathbf{u}_t behaves like a random Gaussian vector. For

$$T \geq 1 + \frac{1}{\varepsilon} \left(\frac{1}{2} - \frac{\log |\beta|}{(k-2) \log n} \right)$$

with probability $1 - O(n^{-c(\log n)^2})$, it holds

$$\widehat{\mathbf{v}} = \mathbf{u}_T = \frac{\tilde{\boldsymbol{\xi}}}{\|\tilde{\boldsymbol{\xi}}\|_2} + O \left(|\beta| \left(\frac{\log n}{\sqrt{n}} \right)^{k-1} \right),$$

where $\tilde{\boldsymbol{\xi}}$ is the standard Gaussian vector in \mathbb{R}^n , the error term is a vector of length bounded by $|\beta|(\log n/\sqrt{n})^{k-1}$.

In Theorem 1, we assume that $\langle \mathbf{u}, \mathbf{v} \rangle \gtrsim 1/\sqrt{n}$, which is generic and is true for a random \mathbf{u} . Moreover, if the initial vector \mathbf{u} is random, then $|\langle \mathbf{u}, \mathbf{v} \rangle| \asymp n^{-1/2}$. Notably, Theorems 1 and 2 together state that power iteration recovers \mathbf{v} if $|\beta| \gg n^{(k-2)/2}$ and fails if $|\beta| \ll n^{(k-2)/2}$. This gives a rigorous proof of the prediction in Richard and Montanari (2014) that the necessary and sufficient condition for the convergence is given by $|\beta| \gtrsim n^{(k-2)/2}$. In practice, it may be possible to use domain knowledge to choose better initialization points. For example, in the classical topic modeling applications Anandkumar et al. (2014b), the unknown vectors \mathbf{v} are related to the topic word distributions, and many documents may be primarily composed of words from just single topic. Therefore, good initialization points can be derived from these single-topic documents.

The special case for $k = 2$, i.e. the spiked matrix model, has been intensively studied since the pioneering work of Johnstone (2001). In this setting it is known [30] that there is an order $O(1)$ critical signal-to-noise ratio, such that below the threshold, it is information-theoretically impossible to recover \mathbf{v} , and above the threshold, the PCA (partially) recovers the unseen eigenvector \mathbf{v} P  ch   (2006); Abbe et al. (2020); O’Rourke et al. (2018); Vu (2011); Zhong (2017); Chen et al. (2018b); Zhang et al. (2018); Cheng et al. (2020). The special case of our results Theorem 1, i.e. the power iteration recovers the eigenvector \mathbf{v} when $\beta \geq n^\varepsilon$, recovers some abovementioned results.

As a consequence of Theorem 1, we have the following central limit theorem for our estimators.

Corollary 3 (Central Limit Theorem) Fix the initialization $\mathbf{u} \in \mathbb{R}^n$ with $\|\mathbf{u}\|_2 = 1$ and $|\langle \mathbf{u}, \mathbf{v} \rangle| \gtrsim 1/\sqrt{n}$. If $|\beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-2}| \geq n^\varepsilon$ with arbitrarily small $\varepsilon > 0$, in Case 1 of Theorem 1, for any fixed unit vector $\mathbf{a} \in \mathbb{R}^n$ obeying

$$|\langle \mathbf{a}, \mathbf{v} \rangle| = o \left(\frac{\beta^3}{\sqrt{n}} \right), \quad (6)$$

and time

$$T \geq 1 + \frac{1}{\varepsilon} \left(\frac{1}{2} + \frac{2 \log |\beta|}{\log n} \right).$$

the estimators $\hat{\mathbf{v}} = \mathbf{u}_T$, and $\hat{\beta} = \mathbf{X}[\hat{\mathbf{v}}^{\otimes k}]$ satisfies

$$\frac{\sqrt{n}\hat{\beta}}{\sqrt{\langle \mathbf{a}, (\mathbf{I}_n - \hat{\mathbf{v}}\hat{\mathbf{v}}^\top) \mathbf{a} \rangle}} \left[\left(1 - \frac{1}{2\hat{\beta}^2}\right)^{-1} \langle \mathbf{a}, \hat{\mathbf{v}} \rangle - \langle \mathbf{a}, \mathbf{v} \rangle \right] \xrightarrow{d} \mathcal{N}(0, 1), \quad (7)$$

as n tends to infinity. We have similar results for Cases 2, 3, 4, by simply changing (β, \mathbf{v}) in (7) to the corresponding limit.

We remark that in Corollary 3, we assume that $|\langle \mathbf{a}, \mathbf{v} \rangle| = o(\beta^3/\sqrt{n})$, which is generic. For example, if \mathbf{v} is delocalized, and \mathbf{a} is supported on finitely many entries, we will have that $|\langle \mathbf{a}, \mathbf{v} \rangle| \lesssim 1/\sqrt{n}$, and (6) is satisfied.

With the central limit theorem for our estimators in Corollary 3, we can easily write down the confidence interval for our estimators.

Corollary 4 (*Prediction Interval*) *Given the asymptotic significance level α , and let $z_\alpha = \Phi(1 - \alpha/2)$ where $\Phi(\cdot)$ is the CDF of a standard Gaussian. If $|\beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-2}| \geq n^\varepsilon$ with arbitrarily small $\varepsilon > 0$, in Case 1 of Theorem 1, for any fixed unit vector $\mathbf{a} \in \mathbb{R}^n$ obeying*

$$|\langle \mathbf{a}, \mathbf{v} \rangle| = o\left(\frac{\beta^3}{\sqrt{n}}\right), \quad (8)$$

and time

$$T \geq 1 + \frac{1}{\varepsilon} \left(\frac{1}{2} + \frac{2 \log |\beta|}{\log n} \right),$$

let $\hat{\mathbf{v}} = \mathbf{u}_T$, and $\hat{\beta} = \mathbf{X}[\hat{\mathbf{v}}^{\otimes k}]$. The asymptotic confidence interval of $\langle \mathbf{a}, \mathbf{v} \rangle$ is given by

$$\frac{1}{1 - 1/(2\hat{\beta}^2)} \left[\langle \mathbf{a}, \hat{\mathbf{v}} \rangle - z_\alpha \frac{\sqrt{\langle \mathbf{a}, (\mathbf{I}_n - \hat{\mathbf{v}}\hat{\mathbf{v}}^\top) \mathbf{a} \rangle}}{\sqrt{n}\hat{\beta}}, \langle \mathbf{a}, \hat{\mathbf{v}} \rangle + z_\alpha \frac{\sqrt{\langle \mathbf{a}, (\mathbf{I}_n - \hat{\mathbf{v}}\hat{\mathbf{v}}^\top) \mathbf{a} \rangle}}{\sqrt{n}\hat{\beta}} \right]. \quad (9)$$

We have similar results for Cases 2, 3, 4, by simply changing (β, \mathbf{v}) in (9) to the corresponding limit.

2.2 General Results: rank- r spiked tensor model

In this section, we state our main results for the general case, the rank- r spiked tensor model (1):

$$\mathbf{X} = \sum_{j=1}^r \beta_j \mathbf{v}_j^{\otimes k} \in \otimes^k \mathbb{R}^n + \mathbf{Z} \quad (10)$$

where

- $\mathbf{X} \in \otimes^k \mathbb{R}^n$ is the k -th order tensor observation.
- $\mathbf{Z} \in \otimes^k \mathbb{R}^n$ is a noise tensor. The entries of \mathbf{Z} are i.i.d. standard $\mathcal{N}(0, 1/n)$ Gaussian random variables.
- $|\beta_1| \geq |\beta_2| \geq \dots \geq |\beta_r|$ are the signal sizes.
- $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r \in \mathbb{R}^n$ are unknown orthonormal vectors to be recovered.

In (10), we assumed that the signals $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$ are orthonormal. The orthogonally decomposable tensor has been widely studied as a benchmark setting for tensor decomposition in the literature Kolda (2001); Chen and Saad (2009); Robeva (2016); Belkin et al. (2018); Auddy and Yuan (2020). Before stating our main results, we need to introduce some more notations and assumptions.

Assumption 5 *We assume that the initialization does not distinguish $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$, such that there exists some large constant $\kappa > 0$*

$$1/\kappa \leq \left| \frac{\langle \mathbf{u}, \mathbf{v}_i \rangle}{\langle \mathbf{u}, \mathbf{v}_j \rangle} \right| \leq \kappa,$$

for all $1 \leq i, j \leq r$.

If we take the uniform initialization, i.e. $\mathbf{u}_0 = \mathbf{u}$ is a uniformly distributed vector in \mathbb{S}^{n-1} . Then with probability $1 - O(r/\sqrt{\kappa})$ we will have $1/\sqrt{\kappa n} \leq |\langle \mathbf{u}, \mathbf{v}_i \rangle| \leq \sqrt{\kappa/n}$ for $1 \leq i \leq r$, and Assumption 5 holds.

The same as in the rank-1 case, the quantities $|\beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle|^{k-2}$ play a crucial role in our power iteration algorithm. We need to make the following technical assumption:

Assumption 6 *Let $j_* = \operatorname{argmax}_j |\beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle|^{k-2}$. We assume that there exists some large constant $\kappa > 0$*

$$(1 - 1/\kappa) |\beta_{j_*} \langle \mathbf{u}, \mathbf{v}_{j_*} \rangle|^{k-2} \geq |\beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle|^{k-2},$$

for all $1 \leq j \leq r$ and $j \neq j_*$.

It turns out under Assumptions 5 and 6, the power iteration converges to \mathbf{v}_{j_*} . Moreover, if we simply take the uniform initialization, i.e. $\mathbf{u}_0 = \mathbf{u}$ is a uniformly distributed vector in \mathbb{S}^{n-1} . Assumption 6 holds for some $1 \leq j_* \leq r$ with probability $1 - O(1/\kappa)$.

Theorem 7 *Fix the initialization $\mathbf{u} \in \mathbb{R}^n$ with $\|\mathbf{u}\|_2 = 1$ and $|\langle \mathbf{u}, \mathbf{v}_j \rangle| \gtrsim 1/\sqrt{n}$, for $1 \leq j \leq r$. Let $j_* = \operatorname{argmax}_j |\beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle|^{k-2}$. Under Assumptions 5 and 6, if $|\beta_{j_*} \langle \mathbf{u}, \mathbf{v}_{j_*} \rangle|^{k-2} \geq n^\varepsilon$ with arbitrarily small $\varepsilon > 0$, then for any fixed unit vector $\mathbf{a} \in \mathbb{R}^n$, and*

$$T \geq 1 + \frac{1}{\varepsilon} \left(\frac{1}{2} + \frac{2 \log |\beta_1|}{\log n} \right) + \frac{\log \log(\sqrt{n} |\beta_1|)}{\log(k-1)},$$

with probability $1 - O(n^{-c(\log n)^2})$, the power iteration estimator $(\hat{\beta}, \hat{\mathbf{v}}) = (\mathbf{X}[\mathbf{u}_T^{\otimes k}], \mathbf{u}_T)$ from (3) satisfies

1. If k is odd, and $\beta_{j_*} > 0$, then $(\widehat{\beta}, \widehat{\mathbf{v}})$ recovers $(\beta_{j_*}, \mathbf{v}_{j_*})$ in the following sense,

$$\begin{aligned} \langle \mathbf{a}, \widehat{\mathbf{v}} \rangle = \langle \mathbf{a}, \mathbf{u}_T \rangle &= \left(1 - \frac{1}{2\beta_{j_*}^2}\right) \langle \mathbf{a}, \mathbf{v}_{j_*} \rangle + \frac{\langle \mathbf{a}, \boldsymbol{\xi} \rangle - \langle \mathbf{a}, \mathbf{v}_{j_*} \rangle \langle \mathbf{v}_{j_*}, \boldsymbol{\xi} \rangle}{\beta_{j_*}} \\ &+ \text{O}_{\mathbb{P}} \left(\frac{\log n}{\sqrt{n}} \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right)^{k-1} + \frac{\log n}{|\beta_1|^2 \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta_1|^{3/2} n^{3/4}} + \frac{1}{|\beta_1|^4} \right), \end{aligned} \quad (11)$$

where $\boldsymbol{\xi} = \mathbf{Z}[\mathbf{v}_{j_*}^{\otimes(k-1)}]$, is an n -dim vector, with each entry i.i.d. $\mathcal{N}(0, 1/n)$ Gaussian random variable. And

$$\begin{aligned} \widehat{\beta} = \mathbf{X}[\mathbf{u}_T^{\otimes k}] &= \beta_{j_*} + \langle \boldsymbol{\xi}, \mathbf{v}_{j_*} \rangle - \frac{k/2 - 1}{\beta_{j_*}} \\ &+ \text{O}_{\mathbb{P}} \left(\frac{\log n}{\sqrt{n}} \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right)^{k-1} + \frac{\log n}{|\beta_1| \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta_1|^{1/2} n^{3/4}} + \frac{1}{|\beta_1|^3} \right). \end{aligned} \quad (12)$$

2. If k is odd, and $\beta_{j_*} < 0$ then $(\widehat{\beta}, \widehat{\mathbf{v}})$ recovers $(-\beta_{j_*}, -\mathbf{v}_{j_*})$, in the sense of (11) and (12) with $(\beta_{j_*}, \mathbf{v}_{j_*})$ replaced by $(-\beta_{j_*}, -\mathbf{v}_{j_*})$
3. If k is even, and $\beta_{j_*} > 0$, then $(\widehat{\beta}, \widehat{\mathbf{v}})$ recovers $(\beta_{j_*}, \text{sgn}(\langle \mathbf{u}, \mathbf{v}_{j_*} \rangle) \mathbf{v}_{j_*})$, in the sense of (11) and (12) with $(\beta_{j_*}, \mathbf{v}_{j_*})$ replaced by $(\beta_{j_*}, \text{sgn}(\langle \mathbf{u}, \mathbf{v}_{j_*} \rangle) \mathbf{v}_{j_*})$;
4. If k is even, and $\beta_{j_*} < 0$, then $(\widehat{\beta}, \widehat{\mathbf{v}})$ recovers $(\beta_{j_*}, (-1)^T \text{sgn}(\langle \mathbf{u}, \mathbf{v}_{j_*} \rangle) \mathbf{v}_{j_*})$, in the sense of (11) and (12) with $(\beta_{j_*}, \mathbf{v}_{j_*})$ replaced by $(\beta_{j_*}, (-1)^T \text{sgn}(\langle \mathbf{u}, \mathbf{v}_{j_*} \rangle) \mathbf{v}_{j_*})$.

In Theorem 7, we assume that $|\langle \mathbf{u}, \mathbf{v}_j \rangle| \gtrsim 1/\sqrt{n}$ for $1 \leq j \leq r$. This is generic and is true for a random initialization \mathbf{u} . Similarly to Theorem 2, if $|\beta_{j_*} \langle \mathbf{u}, \mathbf{v}_{j_*} \rangle^{k-2}| \leq n^{-\varepsilon}$, the iteration does not converge, and \mathbf{u}_t behaves like a random Gaussian vector for any fixed time t .

We want to remark that for multi-rank spiked tensor model, the scenarios for $k = 2$, i.e. the spiked matrix model, and $k \geq 3$ are very different. For the spiked matrix model, in Theorem 7, we always have that $j_* = \text{argmax}_j |\beta_j| = 1$, and power iteration algorithm always converges to the eigenvector corresponding to the largest eigenvalue. However, for rank $k \geq 3$, the power iteration algorithm may converge to any vector \mathbf{v}_j provided that the initialization \mathbf{u} is sufficiently close to \mathbf{v}_j . As a consequence of Theorem 7, we have the following central limit theorem for our estimators.

Corollary 8 Fix the initialization $\mathbf{u} \in \mathbb{R}^n$ with $\|\mathbf{u}\|_2 = 1$ and $|\langle \mathbf{u}, \mathbf{v}_j \rangle| \gtrsim 1/\sqrt{n}$ for $1 \leq j \leq r$. We assume $|\beta \langle \mathbf{u}, \mathbf{v}_{j_*} \rangle^{k-2}| \geq n^\varepsilon$ with arbitrarily small $\varepsilon > 0$, and Assumptions 5 and 6. In Case 1 of Theorem 7, for any fixed unit vector $\mathbf{a} \in \mathbb{R}^n$, for any fixed unit vector $\mathbf{a} \in \mathbb{R}^n$ obeying

$$|\langle \mathbf{a}, \mathbf{v}_{j_*} \rangle| = o\left(\frac{|\beta_1|^3}{\sqrt{n}}\right), \quad (13)$$

and time

$$T \geq 1 + \frac{1}{\varepsilon} \left(\frac{1}{2} + \frac{2 \log |\beta_1|}{\log n} \right),$$

the estimators $\widehat{\mathbf{v}} = \mathbf{u}_T$, and $\widehat{\beta} = \mathbf{X}[\mathbf{u}_T^{\otimes k}]$ satisfy

$$\frac{\sqrt{n}\widehat{\beta}_{j_*}}{\sqrt{\langle \mathbf{a}, (\mathbf{I}_n - \widehat{\mathbf{v}}\widehat{\mathbf{v}}^\top) \mathbf{a} \rangle}} \left[\left(1 - \frac{1}{2\widehat{\beta}_{j_*}^2}\right)^{-1} \langle \mathbf{a}, \widehat{\mathbf{v}} \rangle - \langle \mathbf{a}, \mathbf{v}_{j_*} \rangle \right] \xrightarrow{d} \mathcal{N}(0, 1). \quad (14)$$

We have similar results for Cases 2, 3, 4, by simply changing $(\beta_{j_*}, \mathbf{v}_{j_*})$ in (14) to the corresponding limit.

In the following we take \mathbf{u} to be a random vector uniformly distributed over the unit sphere. The power iteration algorithm can be easily understood in this setting, thanks to Theorem 7. More precisely if $j_* = \operatorname{argmax}_j |\beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle|^{k-2}$ and the initialization \mathbf{u} satisfies Assumptions 5 and 6, then the power iteration estimator $(\widehat{\mathbf{v}}, \widehat{\beta})$ recovers $(\mathbf{v}_{j_*}, \beta_{j_*})$. From the discussions below, for a random vector \mathbf{u} uniformly distributed over the unit sphere, Assumptions 5 and 6 holds with probability $1 - O(1/\sqrt{\kappa})$. We can compute explicitly the probability that index i achieves $\operatorname{argmax}_j |\beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle|^{k-2}$:

$$\begin{aligned} p_i &:= \mathbb{P}(i = \operatorname{argmax}_j |\beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle|^{k-2}) \\ &= \int_0^\infty \sqrt{\frac{2}{\pi}} e^{-x^2/2} \left(\prod_{\ell \neq i} \int_0^{\left(\frac{|\beta_i|}{|\beta_\ell|}\right)^{\frac{1}{k-2}} x} \sqrt{\frac{2}{\pi}} e^{-y^2/2} dy \right) dx, \end{aligned} \quad (15)$$

for any $1 \leq i \leq r$. For spiked matrix model, i.e. $k = 2$, we always have $1 = \operatorname{argmax}_j |\beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle|^{k-2}$, and $p_1 = 1, p_2 = p_3 = \dots = 0$. For spiked tensor models with $k \geq 3$, all those p_i are non-negative and $p_1 \geq p_2 \geq p_3 \geq \dots > 0$. When there exists a large gap between the first signal and the remaining signals, namely $|\beta_1| \geq M|\beta_2|$ for some large constant $M > 0$, then

$$\begin{aligned} p_1 &\geq \int_0^\infty \sqrt{\frac{2}{\pi}} e^{-x^2/2} \left(\prod_{\ell \neq i} \int_0^{M^{\frac{1}{k-2}} x} \sqrt{\frac{2}{\pi}} e^{-y^2/2} dy \right) dx \\ &\geq \int_0^\infty \sqrt{\frac{2}{\pi}} e^{-x^2/2} \left(1 - e^{-M^{\frac{2}{k-2}} x^2/2} \right)^{r-1} dx \\ &\geq \int_0^\infty \sqrt{\frac{2}{\pi}} e^{-x^2/2} \left(1 - (r-1)e^{-M^{\frac{2}{k-2}} x^2/2} \right) dx \geq 1 - \frac{(r-1)}{M^{\frac{2}{k-2}} + 1} = 1 - O\left(M^{-\frac{2}{k-2}}\right), \end{aligned}$$

where in the second line we used the lower bound for the error function. In other words, we can recover the top signal \mathbf{v}_1 with probability $1 - \delta$, provided $|\beta_1|/|\beta_2| \geq C\delta^{-(k-2)/2}$.

Theorem 9 Fix large $\kappa > 0$ and recall p_i as defined (15). If \mathbf{u} is uniformly distributed over the unit sphere, and $|\beta_1| \geq n^{(k-2)/2+\varepsilon}$ with arbitrarily small $\varepsilon > 0$, then for any fixed unit vector $\mathbf{a} \in \mathbb{R}^n$, any $1 \leq i \leq r$, and

$$T \geq 1 + \frac{1}{\varepsilon} \left(\frac{1}{2} + \frac{2 \log |\beta_1|}{\log n} \right) + \frac{\log \log(\sqrt{n}|\beta_1|)}{\log(k-1)},$$

the power iteration estimator $(\widehat{\beta}, \widehat{\mathbf{v}}) = (\mathbf{X}[\mathbf{u}_T^{\otimes k}], \mathbf{u}_T)$ from (3) satisfies

1. If k is odd, and $\beta_i > 0$ then with probability $p_i + O(1/\sqrt{\kappa})$, $(\widehat{\beta}, \widehat{\mathbf{v}})$ recovers (β_i, \mathbf{v}_i) , in the sense

$$\begin{aligned} \langle \mathbf{a}, \widehat{\mathbf{v}} \rangle &= \langle \mathbf{a}, \mathbf{u}_T \rangle = \left(1 - \frac{1}{2\beta_i^2}\right) \langle \mathbf{a}, \mathbf{v}_i \rangle + \frac{\langle \mathbf{a}, \boldsymbol{\xi} \rangle - \langle \mathbf{a}, \mathbf{v}_i \rangle \langle \mathbf{v}_i, \boldsymbol{\xi} \rangle}{\beta_i} \\ &+ O_{\mathbb{P}} \left(\frac{\log n}{\sqrt{n}} \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right)^{k-1} + \frac{\log n}{|\beta_1|^2 \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta_1|^{3/2} n^{3/4}} + \frac{1}{|\beta_1|^4} \right), \end{aligned} \quad (16)$$

where $\boldsymbol{\xi} = \mathbf{Z}[\mathbf{v}_i^{\otimes(k-1)}]$, is an n -dim vector, with each entry i.i.d. $\mathcal{N}(0, 1/n)$ Gaussian random variable. And

$$\begin{aligned} \widehat{\beta} &= \mathbf{X}[\mathbf{u}_T^{\otimes k}] = \beta_i + \langle \boldsymbol{\xi}, \mathbf{v}_i \rangle - \frac{k/2 - 1}{\beta_i} \\ &+ O_{\mathbb{P}} \left(\frac{\log n}{\sqrt{n}} \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right)^{k-1} + \frac{\log n}{|\beta_1| \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta_1|^{1/2} n^{3/4}} + \frac{1}{|\beta_1|^3} \right). \end{aligned} \quad (17)$$

2. If k is odd, and $\beta_i < 0$ then with probability $p_i + O(1/\sqrt{\kappa})$, $(\widehat{\beta}, \widehat{\mathbf{v}})$ recovers $(-\beta_i, -\mathbf{v}_i)$;
 3. If k is even, and $\beta_i > 0$ or $\beta_i < 0$, then with probability $p_i/2 + O(1/\sqrt{\kappa})$, $(\widehat{\beta}, \widehat{\mathbf{v}})$ recovers $(\beta_i, +\mathbf{v}_i)$, and with probability $p_i/2 + O(1/\sqrt{\kappa})$, $(\widehat{\beta}, \widehat{\mathbf{v}})$ recovers $(\beta_i, -\mathbf{v}_i)$.

We want to emphasize here that the scenario for $k = 2$, i.e. the spiked matrix model, and $k \geq 3$ are very different. For spiked matrix model, i.e. $k = 2$, we always have that $p_1 = 0, p_2 = p_3 = \dots = 0$. The power iteration algorithm always converges to the eigenvector corresponding to the largest eigenvalue. We can only recover (β_1, \mathbf{v}_1) no matter how many times we repeat the algorithm. However, for spiked tensor models with $k \geq 3$, all those p_i are nonnegative, $p_1 \geq p_2 \geq p_3 \geq \dots > 0$. By repeating the power iteration algorithm for sufficiently many times, it recovers (β_i, \mathbf{v}_i) with probability roughly p_i .

Similar to the rank one case in Section 2.1, we are also able to establish the asymptotic distribution and confidence interval for multi-rank spiked tensor model with uniformly distributed initialization \mathbf{u} .

Corollary 10 Fix $k \geq 3$, assume \mathbf{u} to be a random vector uniformly distributed over the unit sphere and $|\beta_1| \geq n^{(k-2)/2+\varepsilon}$ with arbitrarily small $\varepsilon > 0$. In Case 1 of Theorem 9, for any fixed unit vector $\mathbf{a} \in \mathbb{R}^n$, and time

$$T \geq 1 + \frac{1}{\varepsilon} \left(\frac{1}{2} + \frac{2 \log |\beta_1|}{\log n} \right) + \frac{\log \log(\sqrt{n}|\beta_1|)}{\log(k-1)},$$

for any $1 \leq i \leq r$, with probability $p_i + O(1/\sqrt{\kappa})$, the estimators $\widehat{\mathbf{v}} = \mathbf{u}_T$ and $\widehat{\beta} = \mathbf{X}[\mathbf{u}_T^{\otimes k}]$ satisfy

$$\frac{\sqrt{n}\widehat{\beta}}{\sqrt{\langle \mathbf{a}, (\mathbf{I}_n - \widehat{\mathbf{v}}\widehat{\mathbf{v}}^\top) \mathbf{a} \rangle}} \left[\langle \mathbf{a}, \widehat{\mathbf{v}} \rangle - \left(1 - \frac{1}{2\widehat{\beta}^2}\right) \langle \mathbf{a}, \mathbf{v}_i \rangle \right] \xrightarrow{d} \mathcal{N}(0, 1). \quad (18)$$

And

$$\sqrt{n} \left(\beta_i - \widehat{\beta} - \frac{k/2 - 1}{\widehat{\beta}} \right) \xrightarrow{d} \mathcal{N}(0, 1). \quad (19)$$

We have similar results for Cases 2, 3, 4, by simply changing (β_i, \mathbf{v}_i) above to the corresponding limit.

We want to emphasize the difference between Corollary 3 and Corollary 10. In the rank one case, the estimators $\hat{\beta}$ and $\langle \mathbf{a}, \hat{\mathbf{v}} \rangle$ are asymptotically Gaussian. In the multi-rank spiked tensor model with $k \geq 3$, those estimators $\hat{\beta}$ and $\langle \mathbf{a}, \hat{\mathbf{v}} \rangle$ are no longer Gaussian. Instead, they are asymptotically a mixture Gaussian with mixture weights $p_1 \geq p_2 \geq p_3 \geq \dots$.

Corollary 11 *Given the asymptotic significance level α , and let $z_\alpha = \Phi(1 - \alpha/2)$ where $\Phi(\cdot)$ is the CDF of a standard Gaussian. Under the conditions in Corollary 10, in Case 1 of Theorem 9, we can find the asymptotic confidence interval of $\langle \mathbf{a}, \mathbf{v}_i \rangle$ as*

$$\frac{1}{1 - 1/(2\hat{\beta}^2)} \left[\langle \mathbf{a}, \hat{\mathbf{v}} \rangle - z_\alpha \frac{\sqrt{\langle \mathbf{a}, (\mathbf{I}_n - \hat{\mathbf{v}}\hat{\mathbf{v}}^\top) \mathbf{a} \rangle}}{\sqrt{n}\hat{\beta}}, \langle \mathbf{a}, \hat{\mathbf{v}} \rangle + z_\alpha \frac{\sqrt{\langle \mathbf{a}, (\mathbf{I}_n - \hat{\mathbf{v}}\hat{\mathbf{v}}^\top) \mathbf{a} \rangle}}{\sqrt{n}\hat{\beta}} \right]$$

and the asymptotic confidence interval of β_i as

$$\left[\hat{\beta} + \frac{k/2 - 1}{\hat{\beta}} - \frac{z_\alpha}{\sqrt{n}}, \hat{\beta} + \frac{k/2 - 1}{\hat{\beta}} + \frac{z_\alpha}{\sqrt{n}} \right].$$

We have similar results for Cases 2, 3, 4, by changing (β_i, \mathbf{v}_i) above to the corresponding limit.

2.3 Proof Sketch

In this section, we outline the proof of Theorem 1. The detailed proofs are given in Appendix A. To analyze the power iteration algorithm (2), we construct an auxiliary iteration, $\mathbf{y}_0 = \mathbf{u}$ and $\mathbf{y}_{t+1} = \mathbf{X}[\mathbf{y}_t^{\otimes(k-1)}]$. This gives us a sequence of vectors $\mathbf{y}_0, \mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_t, \dots$. To analyze \mathbf{y}_{t+1} , we condition on $\mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_t$, then Lemma 12 implies a decomposition of the Gaussian tensor $\mathbf{Z} = \mathbf{Z}_{\text{fixed}} + \mathbf{Z}_{\text{rand}}$ given by a deterministic part and a random part. This can be used to give a decomposition of \mathbf{y}_{t+1}

$$\begin{aligned} \mathbf{y}_{t+1} &= \beta \langle \mathbf{v}, \mathbf{y}_t \rangle^{k-1} \mathbf{v} + \mathbf{Z}_{\text{fixed}}[\mathbf{y}_t^{\otimes(k-1)}] + \mathbf{Z}_{\text{rand}}[\mathbf{y}_t^{\otimes(k-1)}] \\ &=: a_{t+1} \mathbf{v} + b_{t+1} \mathbf{w}_{t+1} + c_{t+1} \boldsymbol{\xi}_{t+1}, \end{aligned} \quad (20)$$

where the first term is the projection in the signal direction and $\boldsymbol{\xi}_{t+1}$ is a random unit vector. The relation (20) induces a recursion for the coefficients $\{a_t, b_t, c_t\}$. Although the recursion is complicated, it can be analyzed using Gaussian concentration estimates. Under the Assumption of Theorem 1 and $|\beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-2}| \geq n^\varepsilon$, we show that eventually a_t will dominate. More precisely, for

$$t \geq 1 + \frac{1}{\varepsilon} \left(\frac{1}{2} + \frac{2 \log |\beta|}{\log n} \right),$$

we have

$$\mathbf{y}_t = a_t \mathbf{v} + \frac{a_t}{\beta} \boldsymbol{\xi} + \text{small error}, \quad (21)$$

where $\boldsymbol{\xi} = \mathbf{Z}[\mathbf{v}^{\otimes(k-1)}]$ behaves like a Gaussian vector. The claims of Theorem 1 then follows from (21).

3. Numerical Study

In this section, we conduct numerical experiments on synthetic data to demonstrate our distributional results provided in Sections 2.1 and 2.2. We fix the dimension $n = 600$ and rank $k = 3$.

3.1 Rank one spiked tensor model

We begin with numerical experiments on rank one case. This section is devoted to numerically studying the efficiency of our estimators for the strength of signals and linear functionals of the signals. We take the signal \mathbf{v} a random vector sampled from the unit sphere in \mathbb{R}^n , and the vector

$$\mathbf{a} = \frac{1}{\sqrt{3}}(\mathbf{e}_{n/3} + \mathbf{e}_{2n/3} + \mathbf{e}_n)$$

For the setting without prior information of the signal, we take the initialization of our power iteration algorithm \mathbf{u} a random vector sampled from the unit sphere in \mathbb{R}^n , and the strength of signal $\beta = n^{(k-2)/2} \approx 24.495$. We plot in Figure 1 our estimators for the strength of signals after normalization

$$\hat{\beta} + \frac{k/2 - 1}{\hat{\beta}} - \beta \tag{22}$$

and our estimators for the linear functionals of the signals

$$\frac{\sqrt{n}\hat{\beta}}{\sqrt{\langle \mathbf{a}, (\mathbf{I}_n - \hat{\mathbf{v}}\hat{\mathbf{v}}^\top) \mathbf{a} \rangle}} \left[\left(1 - \frac{1}{2\hat{\beta}^2}\right)^{-1} \langle \mathbf{a}, \hat{\mathbf{v}} \rangle - \langle \mathbf{a}, \mathbf{v} \rangle \right] \tag{23}$$

as in Corollary 3.

For the setting that there is prior information of the signal, we take the initialization of our power iteration algorithm $\mathbf{u} = (\mathbf{v} + \mathbf{w})/\|\mathbf{v} + \mathbf{w}\|_2$, where \mathbf{v} is a random vector sampled from the unit sphere in \mathbb{R}^n . We plot our estimators for the strength of signals after normalization (22) and our estimators for the linear functionals of the signals (23) for $\beta = 5$ in Figure 2, and for $\beta = 10$ in Figure 3. Although our Theorem 1 and Corollary 3 requires $|\beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-2}| \geq n^\epsilon \gg 1$, Figures 2 and 3 indicate that our estimators $\hat{\beta}$ and $\langle \mathbf{a}, \hat{\mathbf{v}} \rangle$ are asymptotically Gaussian even with small β , i.e. $\beta = 5, 10$. Theorem 1 also indicates that error term in Corollary (3), i.e. the error term in (7), is of order $1/|\beta|$. This matches with our simulation. In Figures 2 and 3, the the difference between the Gaussian fit of our empirical density and the density of $\mathcal{N}(0, 1)$ decreases as β increases from 5 to 10.

In Figure 4, we test the threshold signal-to-noise ratio for the power iteration algorithm. Our Theorems 1 and 2 state that for $|\beta \langle \mathbf{u}_0, \mathbf{v} \rangle^{k-2}| \gg 1$ tensor power iteration recovers the signal \mathbf{v} , and fails when $|\beta \langle \mathbf{u}_0, \mathbf{v} \rangle^{k-2}| \ll 1$. Especially for random initialization, we have that $|\langle \mathbf{u}_0, \mathbf{v} \rangle| \asymp 1/\sqrt{n}$. Our Theorems state that for $|\beta| \gg n^{(k-2)/2}$ tensor power iteration recovers the signal \mathbf{v} , and fails when $|\beta| \ll n^{(k-2)/2}$. Take $k = 3$. In the left panel of Figure 4, we test tensor power iteration with random initialization for various dimensions $n \in \{200, 300, 400, 500, 600\}$ and signal strength $\beta/\sqrt{n} \in (0, 2]$. In the right panel of Figure 4, we test tensor power iteration with fixed small $\beta = 3$ and informative initialization $\beta \langle \mathbf{u}_0, \mathbf{v} \rangle \in (0, 2]$ for various dimensions $n \in \{200, 300, 400, 500, 600\}$. The outputs $\langle \hat{\mathbf{v}}, \mathbf{v} \rangle$ are averaged over 60 independent trials.

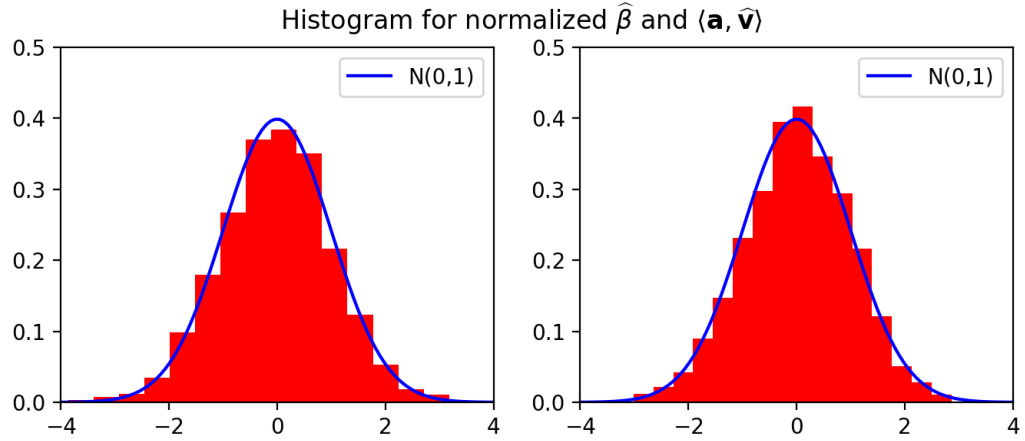


Figure 1: The empirical density of normalized $\hat{\beta}$ as in (22) (left panel), and normalized $\langle \mathbf{a}, \hat{\mathbf{v}} \rangle$ as in (23). The results are reported over 2000 independent trials where the initialization of our power iteration algorithm \mathbf{u} a random vector sampled from the unit sphere in \mathbb{R}^n , and the strength of signal $\beta = n^{(k-2)/2} \approx 24.495$.

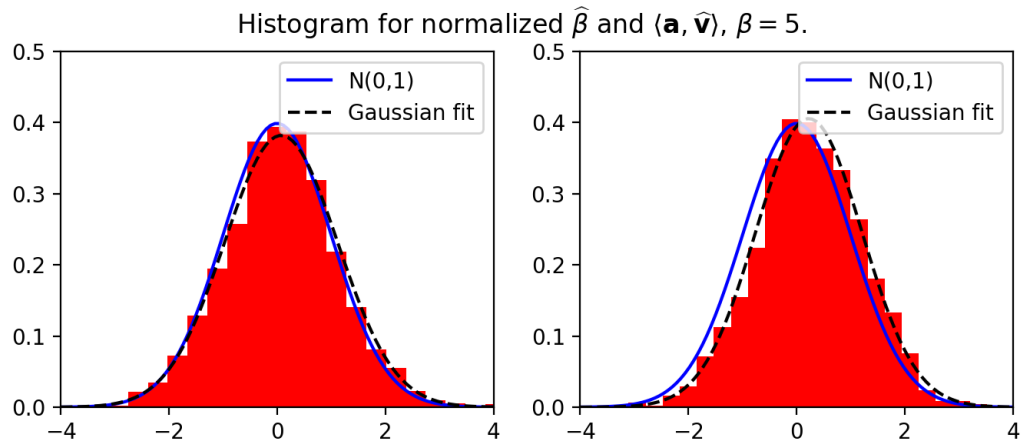


Figure 2: The empirical density of normalized $\hat{\beta}$ as in (22) (left panel), and normalized $\langle \mathbf{a}, \hat{\mathbf{v}} \rangle$ as in (23). The results are reported over 2000 independent trials where the initialization of our power iteration algorithm \mathbf{u} a random vector sampled from the unit sphere in \mathbb{R}^n , and the strength of signal $\beta = 5$.

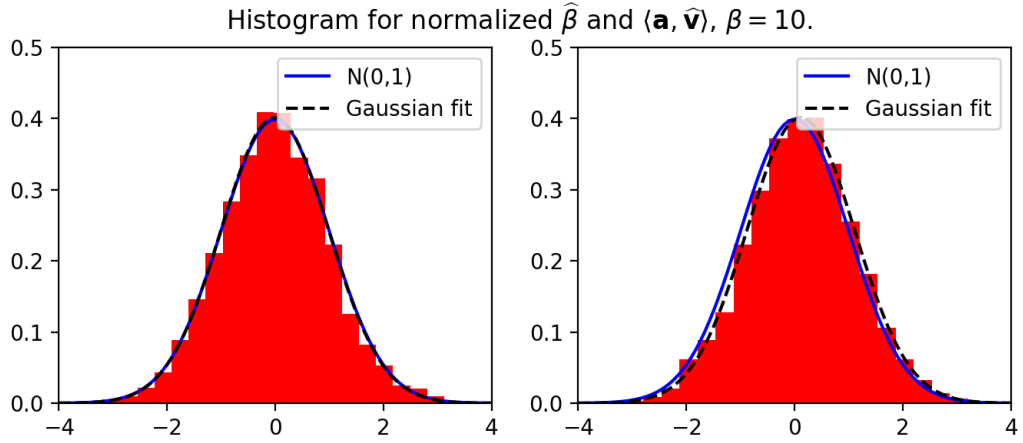


Figure 3: The empirical density of normalized $\hat{\beta}$ as in (22) (left panel), and normalized $\langle \mathbf{a}, \hat{\mathbf{v}} \rangle$ as in (23). The results are reported over 2000 independent trials where the initialization of our power iteration algorithm \mathbf{u} a random vector sampled from the unit sphere in \mathbb{R}^n , and the strength of signal $\beta = 10$.

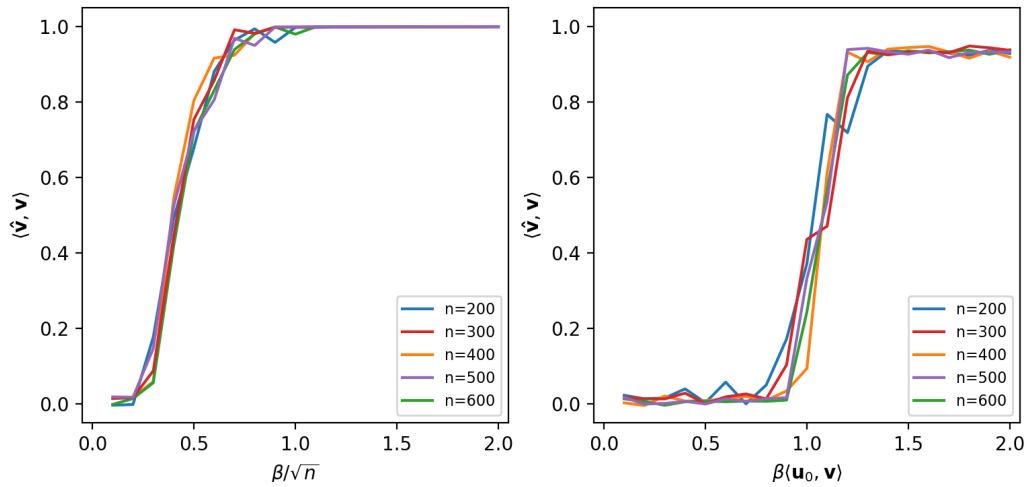


Figure 4: Output of tensor power iteration with random initialization for various signal strength $\beta/\sqrt{n} \in (0, 2]$ (left panel), and tensor power iteration with fixed small $\beta = 3$ and informative initialization $\beta\langle \mathbf{u}_0, \mathbf{v} \rangle \in (0, 2]$.

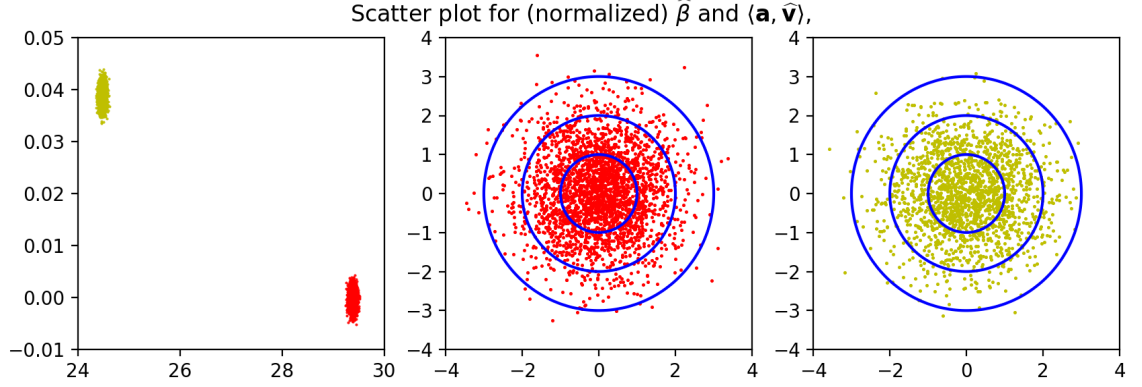


Figure 5: Scatter plot of $(\widehat{\beta}, \langle \mathbf{a}, \widehat{\mathbf{v}} \rangle)$ (first panel), the normalized $(\widehat{\beta}, \langle \mathbf{a}, \widehat{\mathbf{v}} \rangle)$ as in (24) for the cluster corresponding to $(\beta_1, \langle \mathbf{a}, \mathbf{v}_1 \rangle)$ (second panel), the normalized $(\widehat{\beta}, \langle \mathbf{a}, \widehat{\mathbf{v}} \rangle)$ as in (25) for the cluster corresponding to $(\beta_2, \langle \mathbf{a}, \mathbf{v}_2 \rangle)$. The contour plot is a standard 2-dim Gaussian distribution, at 1, 2, 3 standard deviation. The results are reported over 5000 independent trials where the initialization of our power iteration algorithm \mathbf{u} a random vector sampled from the unit sphere in \mathbb{R}^n .

3.2 Rank- r spiked tensor model

In this section, we conduct numerical experiments to demonstrate our distributional results for the multi-rank spiked tensor model. We consider the simplest case that there are two spikes with signals $\mathbf{v}_1, \mathbf{v}_2$, such that they are uniformly sampled from the unit sphere in \mathbb{R}^n and orthogonal to each other $\langle \mathbf{v}_1, \mathbf{v}_2 \rangle = 0$, and the vector

$$\mathbf{a} = \frac{1}{\sqrt{3}}(\mathbf{e}_{n/3} + \mathbf{e}_{2n/3} + \mathbf{e}_n).$$

We test the setting that there is no prior information of the signal. We take the strength of signals $\beta_1 = 1.2 \times n^{(k-2)/2} \approx 29.394$ and $\beta_2 = n^{(k-2)/2} \approx 24.495$ and the initialization of our power iteration algorithm \mathbf{u} a random vector sampled from the unit sphere in \mathbb{R}^n . We scatter plot in Figure 6 our estimator $\widehat{\beta}$ for the strength of signals, and our estimator $\langle \mathbf{a}, \widehat{\mathbf{v}} \rangle$ for the linear functionals of the signals over 5000 independent trials. As seen in the first panel of Figure 6, our estimators $(\widehat{\beta}, \langle \mathbf{a}, \widehat{\mathbf{v}} \rangle)$ form two clusters, centered around $(\beta_1, \langle \mathbf{a}, \mathbf{v}_1 \rangle) \approx (29.394, 0.000)$ and $(\beta_2, \langle \mathbf{a}, \mathbf{v}_2 \rangle) \approx (24.495, 0.039)$. In the second and third panels, we zoom in, and scatter plot for the cluster corresponding to $(\beta_1, \langle \mathbf{a}, \widehat{\mathbf{v}}_1 \rangle) \approx (29.394, 0.000)$

$$\widehat{\beta} + \frac{k/2 - 1}{\widehat{\beta}} - \beta, \quad \frac{\sqrt{n}\widehat{\beta}_1}{\sqrt{\langle \mathbf{a}, (\mathbf{I}_n - \widehat{\mathbf{v}}\widehat{\mathbf{v}}^\top)\mathbf{a} \rangle}} \left[\left(1 - \frac{1}{2\widehat{\beta}^2}\right)^{-1} \langle \mathbf{a}, \widehat{\mathbf{v}} \rangle - \langle \mathbf{a}, \mathbf{v}_1 \rangle \right], \quad (24)$$

and scatter plot for the cluster corresponding to $(\beta_2, \langle \mathbf{a}, \widehat{\mathbf{v}}_2 \rangle) \approx (24.495, 0.039)$

$$\widehat{\beta} + \frac{k/2 - 1}{\widehat{\beta}} - \beta_2, \quad \frac{\sqrt{n}\widehat{\beta}}{\sqrt{\langle \mathbf{a}, (\mathbf{I}_n - \widehat{\mathbf{v}}\widehat{\mathbf{v}}^\top)\mathbf{a} \rangle}} \left[\left(1 - \frac{1}{2\widehat{\beta}^2}\right)^{-1} \langle \mathbf{a}, \widehat{\mathbf{v}} \rangle - \langle \mathbf{a}, \mathbf{v}_2 \rangle \right]. \quad (25)$$

Histogram for normalized $\hat{\beta}$ and $\langle \mathbf{a}, \hat{\mathbf{v}} \rangle$.

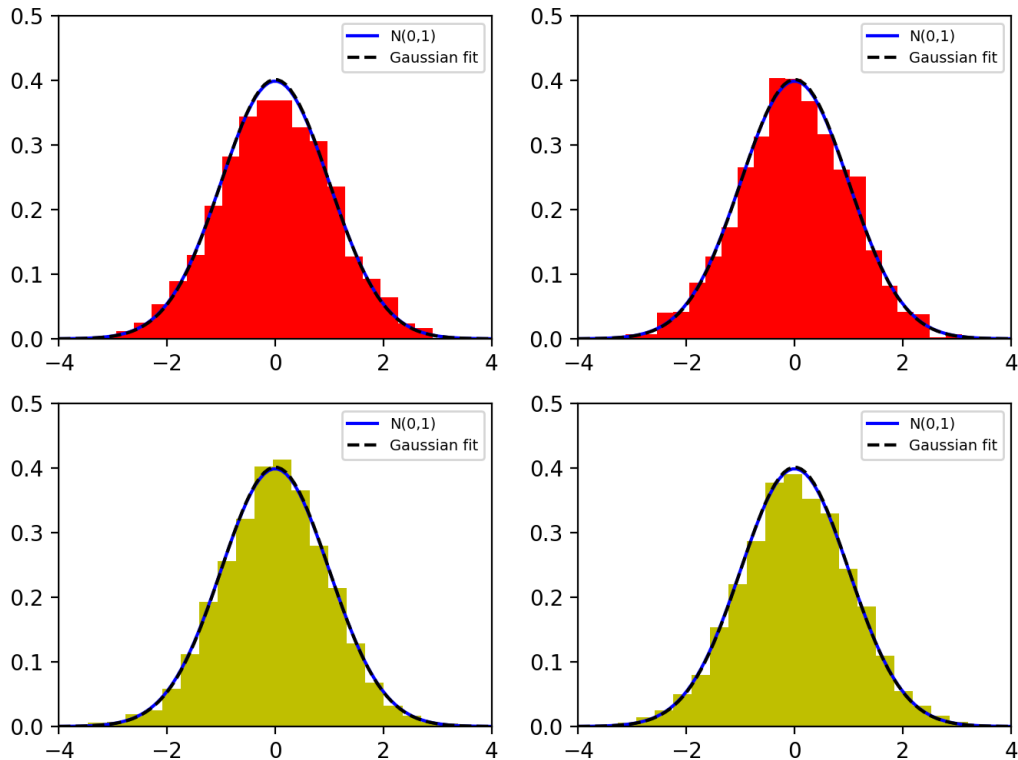


Figure 6: The empirical density of the normalized $(\hat{\beta}, \langle \mathbf{a}, \hat{\mathbf{v}} \rangle)$ as in (24) for the cluster corresponding to $(\beta_1, \langle \mathbf{a}, \mathbf{v}_1 \rangle)$ (second panel), the normalized $(\hat{\beta}, \langle \mathbf{a}, \hat{\mathbf{v}} \rangle)$ as in (25) for the cluster corresponding to $(\beta_2, \langle \mathbf{a}, \mathbf{v}_2 \rangle)$. The results are reported over 5000 independent trials where the initialization of our power iteration algorithm \mathbf{u} a random vector sampled from the unit sphere in \mathbb{R}^n .

As predicted by our Theorem 9, both clusters are asymptotically Gaussian, and the normalized estimators matches pretty well with the contour plot of standard 2-dim Gaussian distribution, at 1, 2, 3 standard deviation.

We plot in Figure 6 our estimators for the strength of signals and the linear functionals of the signals after normalization, for the first cluster (24), and for the second cluster (25).

In Table (1), for each $n \in \{50, 100, 200, 400, 600\}$ and $k = 3$, we take the strength of signals $\beta_1 = 10n^{(k-2)/2}$ and $\beta_2 = 12 \times n^{(k-2)/2}$. Over 1500 independent trials for power iteration with random initialization for each n , we estimate the percentage \hat{p}_1 of estimators converging to β_1 , and the percentage \hat{p}_2 of estimators converging to β_2 . Our theoretical

	$n = 50$	$n = 100$	$n = 200$	$n = 400$	$n = 600$
\widehat{p}_1	0.421	0.421	0.439	0.446	0.445
\widehat{p}_2	0.579	0.579	0.561	0.554	0.555
signal β_1	0.932	0.956	0.948	0.961	0.957
linear form $\langle \mathbf{a}, \mathbf{v}_1 \rangle$	0.965	0.929	0.959	0.945	0.952
signal β_2	0.944	0.944	0.946	0.953	0.949
linear form $\langle \mathbf{a}, \mathbf{v}_2 \rangle$	0.950	0.950	0.937	0.949	0.941

Table 1: Estimated $\widehat{p}_1, \widehat{p}_2$ over 1500 independent trials for dimension $n \in \{50, 100, 200, 400, 600\}$ (top two rows), and numerical coverage rates for our 95% confidence intervals over 1500 independent trials for dimension $n \in \{50, 100, 200, 400, 600\}$ (last four rows).

values are

$$p_1 = \mathbb{P}(|\beta_1 \langle \mathbf{u}, \mathbf{v}_1 \rangle| > |\beta_2 \langle \mathbf{u}, \mathbf{v}_2 \rangle|) \approx 0.44,$$

$$p_2 = \mathbb{P}(|\beta_1 \langle \mathbf{u}, \mathbf{v}_1 \rangle| < |\beta_2 \langle \mathbf{u}, \mathbf{v}_2 \rangle|) \approx 0.56.$$

We also examine the numerical coverage rates for our 95% confidence intervals over 1500 independent trials.

Acknowledgments

The research collaboration was initiated when both G.C. and J.H. were warmly hosted by IAS in the special year of Deep Learning Theory. The research of G.C. is supported by NSF grant NSF-SCALE MoDL (2134209). The research of J.H. is supported by the Simons Foundation as a Junior Fellow at the Simons Society of Fellows, and NSF grant DMS-2054835. The research of Q.Y. is supported by National Natural Science Foundation of China (No. 12101585). We also want to thank the anonymous referees for their helpful comments and suggestions.

Appendix A. Proof of main theorems

A.1 Proof of Theorems 1 and 2

The following lemma on the conditioning of Gaussian tensors will be repeatedly use in the remaining of this section.

Lemma 12 *Let $\mathbf{Z} \in \otimes^k \mathbb{R}^n$ be a random Gaussian tensor. The entries of \mathbf{Z} are i.i.d. standard $\mathcal{N}(0, 1/n)$ Gaussian random variables. Fix $\boldsymbol{\tau}_1, \boldsymbol{\tau}_2, \dots, \boldsymbol{\tau}_t \in \otimes^{k-1} \mathbb{R}^n$ orthonormal $(k-1)$ -th order tensors, i.e. $\langle \boldsymbol{\tau}_i, \boldsymbol{\tau}_j \rangle = \delta_{ij}$, and vectors $\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_t \in \mathbb{R}^n$. Then the distribution of $\mathbf{Z}[\boldsymbol{\tau}]$ conditioned on $\mathbf{Z}[\boldsymbol{\tau}_s] = \boldsymbol{\xi}_s$ for $1 \leq s \leq t$ is*

$$\mathbf{Z}[\boldsymbol{\tau}] \stackrel{d}{=} \sum_{s=1}^t \langle \boldsymbol{\tau}_s, \boldsymbol{\tau} \rangle \boldsymbol{\xi}_s + \tilde{\mathbf{Z}} \left[\boldsymbol{\tau} - \sum_{s=1}^t \langle \boldsymbol{\tau}_s, \boldsymbol{\tau} \rangle \boldsymbol{\tau}_s \right],$$

where $\tilde{\mathbf{Z}}$ is an independent copy of \mathbf{Z} .

Proof [Proof of Lemma 12] For any $(k-1)$ -th order tensor $\boldsymbol{\tau}$, viewed as a vector in $\mathbb{R}^{n^{k-1}}$, we can decompose it as the projection on the span of $\boldsymbol{\tau}_1, \boldsymbol{\tau}_2, \dots, \boldsymbol{\tau}_t$ and the orthogonal part

$$\boldsymbol{\tau} = \sum_{s=1}^t \langle \boldsymbol{\tau}_s, \boldsymbol{\tau} \rangle \boldsymbol{\tau}_s + \left(\boldsymbol{\tau} - \sum_{s=1}^t \langle \boldsymbol{\tau}_s, \boldsymbol{\tau} \rangle \boldsymbol{\tau}_s \right). \quad (26)$$

Using the above decomposition and $\mathbf{Z}[\boldsymbol{\tau}_s] = \boldsymbol{\xi}_s$, we can write $\mathbf{Z}[\boldsymbol{\tau}]$ as

$$\mathbf{Z}[\boldsymbol{\tau}] \stackrel{d}{=} \sum_{s=1}^t \langle \boldsymbol{\tau}_s, \boldsymbol{\tau} \rangle \boldsymbol{\xi}_s + \mathbf{Z} \left[\boldsymbol{\tau} - \sum_{s=1}^t \langle \boldsymbol{\tau}_s, \boldsymbol{\tau} \rangle \boldsymbol{\tau}_s \right], \quad (27)$$

and the first sum and the second term on the righthand side of (27) are independent. The claim (26) follows. \blacksquare

Proof [Proof of Theorem 1] We define an auxiliary iteration, $\mathbf{y}_0 = \mathbf{u}$ and

$$\mathbf{y}_{t+1} = \mathbf{X}[\mathbf{y}_t^{\otimes(k-1)}]. \quad (28)$$

Then with \mathbf{y}_t , our original power iteration (2) is given by $\mathbf{u}_t = \mathbf{y}_t / \|\mathbf{y}_t\|_2$.

Let $\boldsymbol{\xi} = \mathbf{Z}[\mathbf{v}^{\otimes(k-1)}] \in \mathbb{R}^n$. Then the entries of $\boldsymbol{\xi}$ are given by

$$\boldsymbol{\xi}(i) = \mathbf{Z}[\mathbf{v}^{\otimes(k-1)}](i) = \langle \mathbf{Z}, \mathbf{e}_i \otimes \mathbf{v}^{\otimes(k-1)} \rangle = \sum_{i_1, i_2, \dots, i_{k-1} \in [1, n]} \mathbf{Z}_{ii_1 i_2 \dots i_{k-1}} \mathbf{v}(i_1) \mathbf{v}(i_2) \cdots \mathbf{v}(i_{k-1}).$$

From the expression, $\boldsymbol{\xi}(i)$ is a linear combination of Gaussian random variables, itself is also a Gaussian. Moreover, these entries $\boldsymbol{\xi}(i)$ are i.i.d. Gaussian variables with mean zero and variance $1/n$:

$$\mathbb{E}[\boldsymbol{\xi}(i)^2] = \sum_{i_1, i_2, \dots, i_{k-1} \in [1, n]} \mathbb{E}[\mathbf{Z}_{ii_1 i_2 \dots i_{k-1}}^2] \mathbf{v}(i_1)^2 \mathbf{v}(i_2)^2 \cdots \mathbf{v}(i_{k-1})^2 = \frac{1}{n}.$$

We can compute \mathbf{y}_t iteratively: \mathbf{y}_1 is given by

$$\mathbf{y}_1 = \mathbf{X}[\mathbf{y}_0^{\otimes(k-1)}] = \beta \langle \mathbf{y}_0, \mathbf{v} \rangle^{k-1} \mathbf{v} + \mathbf{Z}[\mathbf{y}_0^{\otimes(k-1)}]. \quad (29)$$

For the last term on the righthand side of (29), we can decompose $\mathbf{y}_0^{\otimes(k-1)}$ as a projection on $\mathbf{v}^{\otimes(k-1)}$ and its orthogonal part:

$$\mathbf{y}_0^{\otimes(k-1)} = \langle \mathbf{y}_0, \mathbf{v} \rangle^{k-1} \mathbf{v}^{\otimes(k-1)} + \sqrt{1 - \langle \mathbf{y}_0, \mathbf{v} \rangle^{2(k-1)}} \boldsymbol{\tau}_0,$$

where $\boldsymbol{\tau}_0 \in \otimes^{(k-1)} \mathbb{R}^n$ and $\langle \mathbf{v}^{\otimes(k-1)}, \boldsymbol{\tau}_0 \rangle = 0$, $\langle \boldsymbol{\tau}_0, \boldsymbol{\tau}_0 \rangle = 1$. Thanks to Lemma 12, conditioning on $\boldsymbol{\xi} = \mathbf{Z}[\mathbf{v}^{\otimes(k-1)}]$, $\boldsymbol{\xi}_1 = \mathbf{Z}[\boldsymbol{\tau}_0]$ has the same law as $\tilde{\mathbf{Z}}[\boldsymbol{\tau}_0]$, where $\tilde{\mathbf{Z}}$ is an independent copy of \mathbf{Z} . Since $\langle \boldsymbol{\tau}_0, \boldsymbol{\tau}_0 \rangle = 1$, $\boldsymbol{\xi}_1$ is a Gaussian vector with each entry $\mathcal{N}(0, 1/n)$. With those notations we can rewrite the expression (29) of \mathbf{y}_1 as

$$\mathbf{y}_1 = \beta \langle \mathbf{y}_0, \mathbf{v} \rangle^{k-1} \mathbf{v} + \langle \mathbf{y}_0, \mathbf{v} \rangle^{k-1} \boldsymbol{\xi} + \sqrt{1 - \langle \mathbf{y}_0, \mathbf{v} \rangle^{2(k-1)}} \boldsymbol{\xi}_1. \quad (30)$$

In the following we show that:

Claim 13 *We can compute $\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3, \dots, \mathbf{y}_t$ inductively. The Gram-Schmidt orthonormalization procedure gives an orthogonal base of $\mathbf{v}^{\otimes(k-1)}, \mathbf{y}_0^{\otimes(k-1)}, \mathbf{y}_1^{\otimes(k-1)}, \dots, \mathbf{y}_{t-1}^{\otimes(k-1)}$ as:*

$$\mathbf{v}^{\otimes(k-1)}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}. \quad (31)$$

Let $\boldsymbol{\xi}_{s+1} = \mathbf{Z}[\boldsymbol{\tau}_s]$ for $0 \leq s \leq t-1$. Conditioning on $\boldsymbol{\xi} = \mathbf{Z}[\mathbf{v}^{\otimes(k-1)}]$ and $\boldsymbol{\xi}_{s+1} = \mathbf{Z}[\boldsymbol{\tau}_s]$ for $0 \leq s \leq t-2$, $\boldsymbol{\xi}_t = \mathbf{Z}[\boldsymbol{\tau}_{t-1}]$ is an independent Gaussian vector, with each entry $\mathcal{N}(0, 1/n)$. Then \mathbf{y}_t is in the following form

$$\mathbf{y}_t = a_t \mathbf{v} + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t, \quad b_t \mathbf{w}_t = b_{t0} \boldsymbol{\xi} + b_{t1} \boldsymbol{\xi}_1 + \dots + b_{tt-1} \boldsymbol{\xi}_{t-1}, \quad (32)$$

where $\|\mathbf{w}_t\|_2 = 1$.

Proof [Proof of Claim 13] The Claim 13 for $t = 1$ follows from (30). In the following, assuming Claim 13 holds for t , we prove it for $t+1$.

Let $\mathbf{v}^{\otimes(k-1)}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_t$ be an orthogonal base for $\mathbf{v}^{\otimes(k-1)}, \mathbf{y}_0^{\otimes(k-1)}, \mathbf{y}_1^{\otimes(k-1)}, \dots, \mathbf{y}_t^{\otimes(k-1)}$, obtained by the Gram-Schmidt orthonormalization procedure. More precisely, given those tensors $\mathbf{v}^{\otimes(k-1)}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}$, we denote

$$\begin{aligned} b_{(t+1)0} &= \langle \mathbf{y}_t^{\otimes(k-1)}, \mathbf{v}^{\otimes(k-1)} \rangle, & c_{t+1} &= \langle \mathbf{y}_t^{\otimes(k-1)}, \boldsymbol{\tau}_t \rangle, \\ b_{(t+1)(s+1)} &= \langle \mathbf{y}_t^{\otimes(k-1)}, \boldsymbol{\tau}_s \rangle, & 0 \leq s &\leq t-1. \end{aligned}$$

then $b_{(t+1)0} \mathbf{v}^{\otimes(k-1)} + b_{(t+1)1} \boldsymbol{\tau}_0 + b_{(t+1)2} \boldsymbol{\tau}_1 + \dots + b_{(t+1)t} \boldsymbol{\tau}_{t-1}$ is the projection of $\mathbf{y}_t^{\otimes(k-1)}$ on the span of $\mathbf{v}^{\otimes(k-1)}, \mathbf{y}_0^{\otimes(k-1)}, \mathbf{y}_1^{\otimes(k-1)}, \dots, \mathbf{y}_{t-1}^{\otimes(k-1)}$. With those notations, we can write $\mathbf{y}_t^{\otimes(k-1)}$ as

$$\mathbf{y}_t^{\otimes(k-1)} = b_{(t+1)0} \mathbf{v}^{\otimes(k-1)} + b_{(t+1)1} \boldsymbol{\tau}_0 + b_{(t+1)2} \boldsymbol{\tau}_1 + \dots + b_{(t+1)t} \boldsymbol{\tau}_{t-1} + c_{t+1} \boldsymbol{\tau}_t, \quad (33)$$

Using (32) and (33), we notice that

$$\langle \beta \mathbf{v}^{\otimes k-1}, \mathbf{y}_t^{\otimes(k-1)} \rangle = \beta(a_t + b_t \langle \mathbf{w}_t, \mathbf{v} \rangle + c_t \langle \boldsymbol{\xi}_t, \mathbf{v} \rangle)^{k-1} \mathbf{v},$$

and the iteration (28) implies that

$$\mathbf{y}_{t+1} = \beta(a_t + b_t \langle \mathbf{w}_t, \mathbf{v} \rangle + c_t \langle \boldsymbol{\xi}_t, \mathbf{v} \rangle)^{k-1} \mathbf{v} + b_{t+1} \mathbf{w}_{t+1} + c_{t+1} \mathbf{Z}[\boldsymbol{\tau}_t], \quad (34)$$

where

$$\begin{aligned} b_{t+1} \mathbf{w}_{t+1} &= \mathbf{Z}[b_{(t+1)0} \mathbf{v}^{\otimes(k-1)} + b_{(t+1)1} \boldsymbol{\tau}_0 + b_{(t+1)2} \boldsymbol{\tau}_1 + \cdots + b_{(t+1)t} \boldsymbol{\tau}_{t-1}] \\ &= b_{(t+1)0} \boldsymbol{\xi} + b_{(t+1)1} \boldsymbol{\xi}_1 + b_{(t+1)2} \boldsymbol{\xi}_2 + \cdots + b_{(t+1)t} \boldsymbol{\xi}_t. \end{aligned}$$

Since $\boldsymbol{\tau}_t$ is orthogonal to $\mathbf{v}^{\otimes(k-1)}$, $\boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}$, Lemma 12 implies that conditioning on $\boldsymbol{\xi} = \mathbf{Z}[\mathbf{v}^{\otimes(k-1)}]$ and $\boldsymbol{\xi}_{s+1} = \mathbf{Z}[\boldsymbol{\tau}_s]$ for $0 \leq s \leq t-1$, $\boldsymbol{\xi}_{t+1} = \mathbf{Z}[\boldsymbol{\tau}_t]$ is an independent Gaussian vector, with each entry $\mathcal{N}(0, 1/n)$. The above discussion gives us that

$$\mathbf{y}_{t+1} = a_{t+1} \mathbf{v} + b_{t+1} \mathbf{w}_{t+1} + c_{t+1} \boldsymbol{\xi}_{t+1}, \quad a_{t+1} = \beta(a_t + b_t \langle \mathbf{w}_t, \mathbf{v} \rangle + c_t \langle \boldsymbol{\xi}_t, \mathbf{v} \rangle)^{k-1}. \quad (35)$$

In this way, for any $t \geq 0$, \mathbf{y}_t is given in the form (32). \blacksquare

In the following, We study the case that $\langle \mathbf{u}, \mathbf{v} \rangle > 0$. The case $\langle \mathbf{u}, \mathbf{v} \rangle < 0$ can be proven in exactly the same way, by simply changing (β, \mathbf{v}) with $((-1)^k \beta, -\mathbf{v})$. We prove by induction

Claim 14 *For any fixed time t , with probability at least $1 - O(e^{-c(\log N)^2})$ the following holds: for any $s \leq t$,*

$$\begin{aligned} |a_s| &\gtrsim |\beta| (|b_{s0}| + |b_{s1}| + \cdots + |b_{s(s-1)}|), \\ |a_s| &\gtrsim n^\varepsilon \max\{\mathbf{1}(k \geq 3) |c_s / \beta^{1/(k-2)}|, |c_s / \sqrt{n}|\}. \end{aligned} \quad (36)$$

and

$$\begin{aligned} \|\boldsymbol{\xi}\|, \|\boldsymbol{\xi}_s\|_2 &= 1 + O(\log n / \sqrt{n}), \quad |\langle \mathbf{v}, \boldsymbol{\xi} \rangle|, |\langle \mathbf{a}, \boldsymbol{\xi} \rangle|, |\langle \mathbf{a}, \boldsymbol{\xi}_s \rangle|, \\ \|\text{Proj}_{\text{Span}\{\mathbf{v}, \boldsymbol{\xi}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{s-1}\}}(\boldsymbol{\xi}_s)\|_2 &\lesssim \log n / \sqrt{n}. \end{aligned} \quad (37)$$

Proof [Proof of Claim 14] From (30), $\mathbf{y}_1 = \beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-1} \mathbf{v} + \langle \mathbf{u}, \mathbf{v} \rangle^{k-1} \boldsymbol{\xi} + \sqrt{1 - \langle \mathbf{u}, \mathbf{v} \rangle^{2(k-1)}} \boldsymbol{\xi}_1$. We have $a_1 = \beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-1}$, $b_{10} = \langle \mathbf{u}, \mathbf{v} \rangle^{k-1}$, $b_1 \mathbf{w}_1 = \langle \mathbf{u}, \mathbf{v} \rangle^{k-1} \boldsymbol{\xi}$ and $c_1 = \sqrt{1 - \langle \mathbf{u}, \mathbf{v} \rangle^{2(k-1)}}$. Since $\boldsymbol{\xi}$ is a Gaussian vector with each entry mean zero and variance $1/n$, the concentration for chi-square distribution implies that

$$\|\boldsymbol{\xi}\|_2 = \sqrt{\sum_{i=1}^n \boldsymbol{\xi}(i)^2} = 1 + O(\log n / \sqrt{n})$$

with probability $1 - e^{c(\log n)^2}$. We can directly check that $|a_1| = |\beta b_{10}|$, $|\beta^{1/(k-2)} a_1| = |\beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-2} |^{(k-1)/(k-2)} \gtrsim n^{(k-1)\varepsilon/(k-2)} \geq n^\varepsilon |c_1|$, and $|\sqrt{n} a_1| = |\beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-2}| |\sqrt{n} \langle \mathbf{u}, \mathbf{v} \rangle| \gtrsim n^\varepsilon \geq n^\varepsilon |c_1|$. Moreover, conditioning on $\mathbf{Z}[\mathbf{v}^{\otimes(k-1)}] = \boldsymbol{\xi}$, Lemma 12 implies that $\boldsymbol{\xi}_1 = \mathbf{Z}[\boldsymbol{\tau}_0]$

is an independent Gaussian random vector with each entry $\mathcal{N}(0, 1/n)$. By the standard concentration inequality, it holds that with probability $1 - e^{-c(\log n)^2}$, $\|\boldsymbol{\xi}_1\|_2 = 1 + O(\log n/\sqrt{n})$, $|\langle \mathbf{a}, \boldsymbol{\xi}_1 \rangle|$ and the projection of $\boldsymbol{\xi}_1$ on the span of $\{\mathbf{v}, \boldsymbol{\xi}\}$ is bounded by $\log n/\sqrt{n}$. So far we have proved that (36) and (37) for $t = 1$.

In the following, we assume that (36) holds for t , and prove it for $t + 1$. We recall from (32) and (35) that

$$a_{t+1} = \beta(a_t + b_t \langle \mathbf{w}_t, \mathbf{v} \rangle + c_t \langle \boldsymbol{\xi}_t, \mathbf{v} \rangle)^{k-1}, \quad b_t \mathbf{w}_t = b_{t0} \boldsymbol{\xi} + b_{t1} \boldsymbol{\xi}_1 + \cdots + b_{tt-1} \boldsymbol{\xi}_{t-1} \quad (38)$$

By our induction hypothesis, we have that

$$|b_t \langle \mathbf{w}_t, \mathbf{v} \rangle| \lesssim |b_{t0} \langle \boldsymbol{\xi}, \mathbf{v} \rangle| + |b_{t1} \langle \boldsymbol{\xi}_1, \mathbf{v} \rangle| + \cdots + |b_{t(t-1)} \langle \boldsymbol{\xi}_{t-1}, \mathbf{v} \rangle| \lesssim (\log n/\sqrt{n}) |a_t| / |\beta|, \quad (39)$$

and

$$|c_t \langle \boldsymbol{\xi}_t, \mathbf{v} \rangle| \lesssim (\log n/\sqrt{n}) |c_t| \lesssim (\log n) |a_t| / n^\varepsilon. \quad (40)$$

It follows from plugging (39) and (40) into (38), we get

$$a_{t+1} = \beta(a_t + O(\log n |a_t| / n^\varepsilon))^{k-1} = (1 + O(\log n / n^\varepsilon)) \beta a_t^{k-1}. \quad (41)$$

We recall from (33), the coefficients $b_{(t+1)0}, b_{(t+1)1}, \dots, b_{(t+1)t}$ are determined from the projection of $\mathbf{y}_t^{\otimes(k-1)}$ on $\mathbf{v}^{\otimes(k-1)}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}$

$$\mathbf{y}_t^{\otimes(k-1)} = b_{(t+1)0} \mathbf{v}^{\otimes(k-1)} + b_{(t+1)1} \boldsymbol{\tau}_0 + b_{(t+1)2} \boldsymbol{\tau}_1 + \cdots + b_{(t+1)t} \boldsymbol{\tau}_{t-1} + c_{t+1} \boldsymbol{\tau}_t.$$

We also recall that $\mathbf{v}^{\otimes(k-1)}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}$ are obtained from $\mathbf{v}^{\otimes(k-1)}, \mathbf{y}_0^{\otimes(k-1)}, \mathbf{y}_1^{\otimes(k-1)}, \dots, \mathbf{y}_{t-1}^{\otimes(k-1)}$ by the Gram-Schmidt orthonormalization procedure. So we have that the span of vectors (viewed as vectors) $\mathbf{v}^{\otimes(k-1)}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}$ is the same as the span of the tensors $\mathbf{v}^{\otimes(k-1)}, \mathbf{y}_0^{\otimes(k-1)}, \mathbf{y}_1^{\otimes(k-1)}, \dots, \mathbf{y}_{t-1}^{\otimes(k-1)}$, which is contained in the span of $\{\mathbf{v}, \mathbf{w}_t, \mathbf{y}_0, \dots, \mathbf{y}_{t-1}\}^{\otimes(k-1)}$. Moreover from the relation (32), one can see that the span of $\{\mathbf{v}, \mathbf{w}_t, \mathbf{y}_0, \dots, \mathbf{y}_{t-1}\}$ is the same as the span of $\{\mathbf{v}, \boldsymbol{\xi}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}\}$. It follows that

$$\begin{aligned} & \sqrt{b_{(t+1)0}^2 + b_{(t+1)1}^2 + b_{(t+1)2}^2 + \cdots + b_{(t+1)t}^2} \\ &= \|\text{Proj}_{\text{Span}\{\mathbf{v}^{\otimes(k-1)}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}\}}(a_t \mathbf{v} + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t)^{\otimes(k-1)}\|_2 \\ &\leq \|\text{Proj}_{\text{Span}\{\mathbf{v}, \mathbf{w}_t, \mathbf{y}_0, \dots, \mathbf{y}_{t-1}\}^{\otimes(k-1)}}(a_t \mathbf{v} + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t)^{\otimes(k-1)}\|_2 \\ &\leq \|\text{Proj}_{\text{Span}\{\mathbf{v}, \mathbf{w}_t, \mathbf{y}_0, \dots, \mathbf{y}_{t-1}\}}(a_t \mathbf{v} + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t)\|_2^{k-1} \\ &= \|a_t \mathbf{v} + b_t \mathbf{w}_t + c_t \text{Proj}_{\text{Span}\{\mathbf{v}, \boldsymbol{\xi}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}\}}(\boldsymbol{\xi}_t)\|_2^{k-1} \\ &\lesssim \left(|a_t| + |b_t| + \frac{\log n |c_t|}{\sqrt{n}} \right)^{k-1} \lesssim |a_t|^{k-1} \lesssim |a_{t+1}| / |\beta|, \end{aligned} \quad (42)$$

where in the last line we used our induction hypothesis that $\|\text{Proj}_{\text{Span}\{\mathbf{v}, \boldsymbol{\xi}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}\}}(\boldsymbol{\xi}_t)\|_2 \lesssim \log n / \sqrt{n}$.

Finally we estimate c_{t+1} . We recall from (33), the coefficient c_{t+1} is the remainder of $\mathbf{y}_t^{\otimes(k-1)}$ after projecting on $\mathbf{v}^{\otimes(k-1)}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}$. It is bounded by the remainder of $\mathbf{y}_t^{\otimes(k-1)}$ after projecting on $\mathbf{v}^{\otimes(k-1)}$,

$$|c_{t+1}| \leq \|\mathbf{y}_t^{\otimes(k-1)} - a_t^{k-1} \mathbf{v}^{\otimes(k-1)}\|_2 = \|(a_t \mathbf{v} + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t)^{\otimes(k-1)} - a_t^{k-1} \mathbf{v}^{\otimes(k-1)}\|_2.$$

The difference $(a_t \mathbf{v} + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t)^{\otimes(k-1)} - a_t^{k-1} \mathbf{v}^{\otimes(k-1)}$ is a sum of terms in the following form,

$$\boldsymbol{\eta}_1 \otimes \boldsymbol{\eta}_2 \otimes \dots \otimes \boldsymbol{\eta}_{k-1}, \quad (43)$$

where vectors $\boldsymbol{\eta}_1, \boldsymbol{\eta}_2, \dots, \boldsymbol{\eta}_{k-1} \in \{a_t \mathbf{v}, b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t\}$, and at least one of them is $b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t$. We notice that by our induction hypothesis, $\|b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t\|_2 \lesssim |b_t| \|\mathbf{w}_t\|_2 + |c_t| \|\boldsymbol{\xi}_t\|_2 \lesssim |b_t| + |c_t|$. For the L_2 norm of (43), each copy of $a_t \mathbf{v}$ contributes a_t and each copy of $b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t$ contributes a factor $|b_t| + |c_t|$. We conclude that

$$|c_{t+1}| \leq \|(a_t \mathbf{v} + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t)^{\otimes(k-1)} - a_t^{k-1} \mathbf{v}^{\otimes(k-1)}\|_2 \lesssim \sum_{r=1}^{k-1} |a_t|^{k-1-r} (|b_t| + |c_t|)^r. \quad (44)$$

Combining the above estimate with (41) that $|a_{t+1}| \asymp |\beta| |a_t|^{k-1}$, we divide both sides of (44) by $|\beta| |a_t|^{k-1}$,

$$\frac{|c_{t+1}|}{|a_{t+1}|} \lesssim \frac{1}{|\beta|} \sum_{r=1}^{k-1} \left(\frac{|b_t|}{|a_t|} + \frac{|c_t|}{|a_t|} \right)^r \lesssim \frac{1}{|\beta|} \sum_{r=1}^{k-1} \left(\frac{1}{|\beta|} + \frac{|c_t|}{|a_t|} \right)^r, \quad (45)$$

where we used our induction hypothesis that $|a_t| \gtrsim |\beta| |b_t|$. There are three cases:

1. If $|c_t|/|a_t| \geq 1$, then

$$\frac{|c_{t+1}|}{|a_{t+1}|} \lesssim \frac{1}{|\beta|} \sum_{r=1}^{k-1} \left(\frac{1}{|\beta|} + \frac{|c_t|}{|a_t|} \right)^r \lesssim \frac{1}{|\beta|} \left(\frac{|c_t|}{|a_t|} \right)^{k-1}.$$

If $k = 2$, then our assumption $|\beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-2}| = |\beta| \geq n^\varepsilon$, implies that $|c_{t+1}|/|a_{t+1}| \lesssim (|c_t|/|a_t|)/n^\varepsilon$. If $k \geq 2$, by our induction hypothesis $|c_t|/|a_t| \lesssim \beta^{1/(k-2)}/n^\varepsilon$. This implies $(|c_t|/|a_t|)^{k-2}/|\beta| \lesssim 1/n^\varepsilon$, and we still get that $|c_{t+1}|/|a_{t+1}| \lesssim (|c_t|/|a_t|)/n^\varepsilon$.

2. If $1/|\beta| \lesssim |c_t|/|a_t| \leq 1$, then

$$\frac{|c_{t+1}|}{|a_{t+1}|} \lesssim \frac{1}{|\beta|} \sum_{r=1}^{k-1} \left(\frac{1}{|\beta|} + \frac{|c_t|}{|a_t|} \right)^r \lesssim \frac{1}{|\beta|} \left(\frac{|c_t|}{|a_t|} \right) \lesssim \frac{1}{n^\varepsilon} \left(\frac{|c_t|}{|a_t|} \right),$$

where we used that $|\beta| \geq |\beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-2}| \geq n^\varepsilon$.

3. Finally for $|c_t|/|a_t| \lesssim 1/|\beta|$, we will have

$$\frac{|c_{t+1}|}{|a_{t+1}|} \lesssim \frac{1}{|\beta|} \sum_{r=1}^{k-1} \left(\frac{1}{|\beta|} + \frac{|c_t|}{|a_t|} \right)^r \lesssim \frac{1}{|\beta|} \left(\frac{1}{|\beta|} \right) \lesssim \frac{1}{|\beta|^2}.$$

In all these cases if $|c_t|/|a_t| \lesssim \min\{\sqrt{n}, \mathbf{1}(k \geq 3)|\beta|^{1/(k-2)}\}/n^\varepsilon$, we have $|c_{t+1}|/|a_{t+1}| \lesssim \min\{\sqrt{n}, \mathbf{1}(k \geq 3)|\beta|^{1/(k-2)}\}/n^\varepsilon$. This finishes the proof of the induction (36).

For (37), since $\boldsymbol{\tau}_t$ is orthogonal to $\mathbf{v}^{\otimes(k-1)}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}$, Lemma 12 implies that conditioning on $\boldsymbol{\xi} = \mathbf{Z}[\mathbf{v}^{\otimes(k-1)}]$ and $\boldsymbol{\xi}_{s+1} = \mathbf{Z}[\boldsymbol{\tau}_s]$ for $0 \leq s \leq t-1$, $\boldsymbol{\xi}_{t+1} = \mathbf{Z}[\boldsymbol{\tau}_t]$ is an independent Gaussian vector, with each entry $\mathcal{N}(0, 1/n)$. By the standard concentration inequality, it holds that with probability $1 - e^{c(\log n)^2}$, $\|\boldsymbol{\xi}_{t+1}\|_2 = 1 + O(\log n/\sqrt{n})$, $|\langle \mathbf{a}, \boldsymbol{\xi}_{t+1} \rangle|$ and the projection of $\boldsymbol{\xi}_{t+1}$ on the span of $\{\mathbf{v}, \boldsymbol{\xi}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_t\}$ is bounded by $\log n/\sqrt{n}$. This finishes the proof of the induction (37). \blacksquare

Next, using (36) and (37) in Claim 14 as input, we prove that for

$$t \geq 1 + \frac{1}{\varepsilon} \left(\frac{1}{2} + \frac{2 \log |\beta|}{\log n} \right), \quad (46)$$

with probability $1 - e^{c(\log n)^2}$ we have

$$\mathbf{y}_t = a_t \mathbf{v} + b_{t0} \boldsymbol{\xi} + b_{t1} \boldsymbol{\xi}_1 + \dots + b_{tt-1} \boldsymbol{\xi}_{t-1} + c_t \boldsymbol{\xi}_t, \quad (47)$$

such that

$$b_{t0} = \frac{a_t}{\beta} + O\left(\frac{\log n |a_t|}{|\beta|^2 \sqrt{n}}\right) \quad |b_{t1}|, |b_{t2}|, \dots, |b_{tt-1}| \lesssim \frac{(\log n)^{1/2} |a_t|}{|\beta|^{3/2} n^{1/4}}, \quad |c_t| \lesssim |a_t|/\beta^2. \quad (48)$$

Let $x_t = |c_t/a_t| \ll |\beta|^{1/(k-2)}$, then (45) implies

$$x_{t+1} \lesssim \frac{1}{|\beta|} \sum_{r=1}^{k-1} \left(\frac{1}{|\beta|} + x_t \right)^r,$$

from the discussion after (45), we have that either $x_{t+1} \lesssim 1/\beta^2$, or $x_{t+1} \lesssim x_t/n^\varepsilon$. Since $x_1 = |c_1/a_1| \lesssim n^{1/2-\varepsilon}$, we conclude that it holds

$$x_t = |c_t/a_t| \lesssim 1/\beta^2, \quad \text{when } t \geq \frac{1}{\varepsilon} \left(\frac{1}{2} + \frac{2 \log |\beta|}{\log n} \right). \quad (49)$$

To derive the upper bound of $b_{t1}, b_{t2}, \dots, b_{tt-1}$, we use (42).

$$\begin{aligned} & b_{(t+1)0}^2 + b_{(t+1)1}^2 + b_{(t+1)2}^2 + \dots + b_{(t+1)t}^2 \\ & \leq \|a_t \mathbf{v} + b_t \mathbf{w}_t + c_t \text{Proj}_{\text{Span}\{\mathbf{v}, \boldsymbol{\xi}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}\}}(\boldsymbol{\xi}_t)\|_2^{2(k-1)} \\ & = \left(a_t^2 + O\left(|a_t| (|b_t| + |c_t|) \frac{\log n}{\sqrt{n}} + \left(|b_t| + |c_t| \frac{\log n}{\sqrt{n}} \right)^2 \right) \right)^{k-1}, \end{aligned} \quad (50)$$

where we used our induction (37) that $|\langle \boldsymbol{\xi}, \mathbf{v} \rangle|, |\langle \boldsymbol{\xi}_1, \mathbf{v} \rangle|, \dots, |\langle \boldsymbol{\xi}_t, \mathbf{v} \rangle| \lesssim \log n/\sqrt{n}$ and the projection $\|\text{Proj}_{\text{Span}\{\mathbf{v}, \boldsymbol{\xi}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}\}}(\boldsymbol{\xi}_t)\|_2 \lesssim \log n/\sqrt{n}$. Moreover, the first term $b_{(t+1)0}$ is the projection of $\mathbf{y}_t^{\otimes(k-1)}$ on $\mathbf{v}^{\otimes(k-1)}$,

$$b_{(t+1)0} = \langle a_t \mathbf{v} + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t, \mathbf{v} \rangle^{k-1} = \left(a_t + O\left(\frac{\log n (|b_t| + |c_t|)}{\sqrt{n}} \right) \right)^{k-1}, \quad (51)$$

where we used (37) that $|\langle \boldsymbol{\xi}, \mathbf{v} \rangle|, |\langle \boldsymbol{\xi}_1, \mathbf{v} \rangle|, \dots, |\langle \boldsymbol{\xi}_t, \mathbf{v} \rangle| \lesssim \log n / \sqrt{n}$. Now we can take difference of (50) and (51), and use that $|b_t| \lesssim |a_t|/|\beta|$ from (36) and $|c_t| \lesssim |a_t|/|\beta|$ from (49),

$$b_{(t+1)0} = a_t^{k-1} + \mathcal{O}\left(|a_t|^{k-1} \frac{\log n}{|\beta| \sqrt{n}}\right), \quad b_{(t+1)1}^2 + b_{(t+1)2}^2 + \dots + b_{(t+1)t}^2 \lesssim a_t^{2(k-1)} \frac{\log n}{|\beta| \sqrt{n}}. \quad (52)$$

From (38) and (41), we have that

$$a_{t+1} = \beta b_{(t+1)0} \asymp \beta a_t^{k-1}. \quad (53)$$

Using the above relation, we can simplify (52) as

$$b_{(t+1)0} = \frac{a_{t+1}}{\beta} + \mathcal{O}\left(\frac{\log n |a_{t+1}|}{|\beta|^2 \sqrt{n}}\right), \quad |b_{(t+1)1}|, |b_{(t+1)2}|, \dots, |b_{(t+1)t}| \lesssim \frac{(\log n)^{1/2} |a_{t+1}|}{|\beta|^{3/2} n^{1/4}}.$$

This finishes the proof of (48).

With the expression (48), we can process to prove our main results (4) and (5). Thanks to (37), for t satisfies (46), we have that with probability at least $1 - \mathcal{O}(e^{-c(\log N)^2})$

$$\|\mathbf{y}_t\|_2^2 = a_t^2 \left(1 + \frac{1}{\beta^2} + \frac{2\langle \mathbf{v}, \boldsymbol{\xi} \rangle}{\beta} + \mathcal{O}\left(\frac{\log n}{\beta^2 \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta|^{3/2} n^{3/4}} + \frac{1}{\beta^4}\right)\right). \quad (54)$$

By rearranging it we get

$$1/\|\mathbf{y}_t\|_2 = \frac{1}{|a_t|} \left(1 - \frac{1}{2\beta^2} - \frac{\langle \mathbf{v}, \boldsymbol{\xi} \rangle}{\beta} + \mathcal{O}\left(\frac{\log n}{\beta^2 \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta|^{3/2} n^{3/4}} + \frac{1}{\beta^4}\right)\right). \quad (55)$$

We can take the inner product $\langle \mathbf{a}, \mathbf{y}_t \rangle$ using (47) and (48), and multiply (55)

$$\begin{aligned} \langle \mathbf{a}, \mathbf{u}_t \rangle &= \frac{\langle \mathbf{a}, \mathbf{y}_t \rangle}{\|\mathbf{y}_t\|_2} = \text{sgn}(a_t) \left(\left(1 - \frac{1}{2\beta^2}\right) \langle \mathbf{a}, \mathbf{v} \rangle + \frac{\langle \mathbf{a}, \boldsymbol{\xi} \rangle - \langle \mathbf{a}, \mathbf{v} \rangle \langle \mathbf{v}, \boldsymbol{\xi} \rangle}{\beta} \right) \\ &\quad + \mathcal{O}_{\mathbb{P}}\left(\frac{\log n}{\beta^2 \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta|^{3/2} n^{3/4}} + \frac{|\langle \mathbf{a}, \mathbf{v} \rangle|}{\beta^4}\right), \end{aligned} \quad (56)$$

where we used (37) that with high probability $|\langle \mathbf{a}, \boldsymbol{\xi}_s \rangle|, |\langle \mathbf{a}, \boldsymbol{\xi}_s \rangle|$ for $1 \leq s \leq t$ are bounded by $\log n / \sqrt{n}$. This finishes the proof of (4). For $\hat{\beta}$ in (5), we have

$$\mathbf{X}[\mathbf{u}_t^{\otimes k}] = \frac{\mathbf{X}[\mathbf{y}_t^{\otimes k}]}{\|\mathbf{y}_t\|_2^k} = \frac{\langle \mathbf{y}_t, \mathbf{X}[\mathbf{y}_t^{\otimes(k-1)}] \rangle}{\|\mathbf{y}_t\|_2^k} = \frac{\langle \mathbf{y}_t, \mathbf{y}_{t+1} \rangle}{\|\mathbf{y}_t\|_2^k}. \quad (57)$$

Thanks to (49), (35) and (37), for t satisfies (46), with probability at least $1 - \mathcal{O}(e^{-c(\log N)^2})$, we have

$$\mathbf{y}_{t+1} = a_{t+1} \mathbf{v} + b_{(t+1)0} \boldsymbol{\xi} + b_{(t+1)1} \boldsymbol{\xi}_1 + \dots + b_{(t+1)t} \boldsymbol{\xi}_t + c_{t+1} \boldsymbol{\xi}_{t+1},$$

where $|c_{t+1}| \lesssim |a_t|^{k-1}/\beta^2$,

$$\begin{aligned} a_{t+1} &= \beta(a_t + b_{t0}\langle \boldsymbol{\xi}, \mathbf{v} \rangle + b_{t1}\langle \boldsymbol{\xi}_1, \mathbf{v} \rangle + \cdots + b_{t(t-1)}\langle \boldsymbol{\xi}_{t-1}, \mathbf{v} \rangle + c_t\langle \boldsymbol{\xi}_t, \mathbf{v} \rangle)^{k-1} \\ &= \beta a_t^{k-1} \left(1 + \frac{\langle \boldsymbol{\xi}, \mathbf{v} \rangle}{\beta} + \mathcal{O} \left(\frac{\log n}{\beta^2 \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta|^{3/2} n^{3/4}} \right) \right)^{k-1}, \end{aligned} \quad (58)$$

and

$$b_{(t+1)0} = a_t^{k-1} \left(1 + \mathcal{O} \left(\frac{\log n}{|\beta| \sqrt{n}} \right) \right), \quad |b_{(t+1)1}|, |b_{(t+1)2}| + \cdots + |b_{(t+1)t}| \lesssim a_t^{k-1} \frac{(\log n)^{1/2}}{|\beta|^{1/2} n^{1/4}}.$$

From the discussion above, combining with (47) and (48) with straightforward computation, we have

$$\langle \mathbf{y}_t, \mathbf{y}_{t+1} \rangle = \beta a_t^k \left(1 + \frac{1}{\beta^2} + \frac{(k+1)\langle \boldsymbol{\xi}, \mathbf{v} \rangle}{\beta} + \mathcal{O} \left(\frac{\log n}{\beta^2 \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta|^{3/2} n^{3/4}} \right) \right). \quad (59)$$

By plugging (55) and (59) into (57), we get

$$\mathbf{X}[\mathbf{u}_t^{\otimes k}] = \text{sgn}(a_t)^k \left(\beta + \langle \boldsymbol{\xi}, \mathbf{v} \rangle - \frac{k/2 - 1}{\beta} \right) + \mathcal{O} \left(\frac{\log n}{|\beta| \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta|^{1/2} n^{3/4}} + \frac{1}{|\beta|^3} \right)$$

Since by our assumption, in Case 1 we have that $\beta > 0$. Thanks to (58) $a_{t+1} = \beta a_t^{k-1}(1 + o(1))$, especially a_{t+1} and a_t are of the same sign. In the case $\langle \mathbf{u}, \mathbf{v} \rangle > 0$, we have $a_1 = \beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-1} > 0$. We conclude that $a_t > 0$. Therefore $\text{sgn}(\mathbf{X}[\mathbf{u}_t^{\otimes k}]) = \text{sgn}(a_t)^k = +$, and it follows that

$$\mathbf{X}[\mathbf{u}_t^{\otimes k}] = \beta + \langle \boldsymbol{\xi}, \mathbf{v} \rangle - \frac{k/2 - 1}{\beta} + \mathcal{O} \left(\frac{\log n}{|\beta| \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta|^{1/2} n^{3/4}} + \frac{1}{|\beta|^3} \right)$$

This finishes the proof of (5). The Cases 2, 3, 4 follow by simply changing (β, \mathbf{v}) in the righthand side of (4) and (5) to the corresponding limit. \blacksquare

Proof [Proof of Theorem 2]

We use the same notations as in the proof of Theorem 1. If $|\beta| \geq n^\varepsilon$ and $|\beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-2}| \leq n^{-\varepsilon}$, then we first prove by induction that for any fixed time t , with probability at least $1 - \mathcal{O}(e^{-c(\log N)^2})$ the following holds: for any $s \leq t$,

$$\begin{aligned} |b_{s0}|, |b_{s1}|, \dots, |b_{s(s-1)}| &\lesssim \max\{|c_s|/|\beta|^{(k-1)/(k-2)}, (\log n)^{k-1} |c_s|/n^{(k-1)/2}\}, \\ |c_s| &\geq n^\varepsilon \beta^{1/(k-2)} |a_s|, \end{aligned} \quad (60)$$

and

$$\begin{aligned} \|\boldsymbol{\xi}\|, \|\boldsymbol{\xi}_s\|_2 &= 1 + \mathcal{O}(\log n / \sqrt{n}), \\ |\langle \mathbf{v}, \boldsymbol{\xi} \rangle|, \|\text{Proj}_{\text{Span}\{\mathbf{v}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{s-1}\}}(\boldsymbol{\xi}_s)\|_2 &\lesssim \log n / \sqrt{n}. \end{aligned} \quad (61)$$

From (30), $a_1 = \beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-1}$, $b_{10} = \langle \mathbf{u}, \mathbf{v} \rangle^{k-1}$ and $c_1 = \sqrt{1 - \langle \mathbf{u}, \mathbf{v} \rangle^{2(k-1)}}$. Since $|\beta| \geq n^\varepsilon$ and $|\beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-2}| \leq n^{-\varepsilon}$, we have that $|\langle \mathbf{u}, \mathbf{v} \rangle| \leq n^{-2\varepsilon/(k-2)} \ll 1$ and therefore $|c_1| \asymp 1$. We can check that $|\beta^{1/(k-2)} a_1| = |\beta \langle \mathbf{u}, \mathbf{v} \rangle^{k-2}|^{(k-1)/(k-2)} \leq n^{-\varepsilon} \lesssim n^{-\varepsilon} |c_1|$ and $|b_{10}| = |a_1/\beta| \lesssim n^{-\varepsilon} |c_1/\beta^{(k-1)/(k-2)}|$. Moreover, conditioning on $\mathbf{Z}[\mathbf{v}^{\otimes(k-1)}] = \boldsymbol{\xi}$, Lemma 12 implies that $\boldsymbol{\xi}_1 = \mathbf{Z}[\boldsymbol{\tau}_0]$ is an independent Gaussian random vector with each entry $\mathcal{N}(0, 1/n)$. By the standard concentration inequality, it holds that with probability $1 - e^{c(\log n)^2}$, $\|\boldsymbol{\xi}_1\|_2 = 1 + O(\log n/\sqrt{n})$, and the projection of $\boldsymbol{\xi}_1$ on the span of $\{\mathbf{v}, \boldsymbol{\xi}\}$ is bounded by $\log n/\sqrt{n}$. So far we have proved (60) and (61) for $t = 1$.

In the following, assuming the statements (60) and (61) hold for t , we prove them for $t + 1$. From (38), using (38) and (39), we have

$$\begin{aligned} |a_{t+1}| &= \left| \beta(a_t + b_t \langle \mathbf{w}_t, \mathbf{v} \rangle + c_t \langle \boldsymbol{\xi}_t, \mathbf{v} \rangle)^{k-1} \right| \\ &\lesssim |\beta| \left(|a_t| + \frac{\log n(|b_{t0}| + |b_{t1}| + \dots + |b_{t(t-1)}|)}{\sqrt{n}} + \frac{\log n |c_t|}{\sqrt{n}} \right)^{k-1} \\ &\lesssim |\beta| \left(|a_t| + \frac{\log n |c_t|}{\sqrt{n}} \right)^{k-1} \lesssim |\beta| \left(\frac{|c_t|}{n^\varepsilon |\beta|^{1/(k-2)}} + \frac{\log n |c_t|}{\sqrt{n}} \right)^{k-1} \\ &\lesssim |\beta| |c_t|^{k-1} \left(\frac{1}{n^\varepsilon |\beta|^{1/(k-2)}} + \frac{\log n}{\sqrt{n}} \right)^{k-1} \lesssim \frac{|c_t|^{k-1}}{n^\varepsilon |\beta|^{1/(k-2)}}, \end{aligned} \quad (62)$$

where in the third line we used our induction hypothesis that $|b_{t0}| + |b_{t1}| + \dots + |b_{t(t-1)}| \lesssim |c_t|$, and $n^{-\varepsilon} \geq |\beta| |\langle \mathbf{u}, \mathbf{v} \rangle|^{k-2} \gtrsim |\beta|/n^{(k-2)/2}$.

For $b_{(t+1)0}, b_{(t+1)1}, \dots, b_{(t+1)t}$, from (42) we have

$$\begin{aligned} \sqrt{b_{(t+1)0}^2 + b_{(t+1)1}^2 + b_{(t+1)2}^2 + \dots + b_{(t+1)t}^2} &\lesssim \left(|a_t| + |b_t| + \frac{\log n |c_t|}{\sqrt{n}} \right)^{k-1} \\ &\lesssim \left(|a_t| + \frac{\log n |c_t|}{\sqrt{n}} \right)^{k-1} \lesssim \left(\frac{|c_t|}{n^\varepsilon |\beta|^{1/(k-2)}} + \frac{\log n |c_t|}{\sqrt{n}} \right)^{k-1} \\ &\lesssim |c_t|^{k-1} \left(\frac{1}{n^\varepsilon |\beta|^{1/(k-2)}} + \frac{\log n}{\sqrt{n}} \right)^{k-1}. \end{aligned} \quad (63)$$

Finally we estimate c_{t+1} . We recall from (33), the coefficient c_{t+1} is the remainder of $\mathbf{y}_t^{\otimes(k-1)}$ after projecting on $\mathbf{v}^{\otimes(k-1)}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}$. We have the following lower bound for c_{t+1}

$$\begin{aligned} |c_{t+1}|^2 &= \|(a_t \mathbf{v} + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t)^{\otimes(k-1)}\|_2^2 - (b_{(t+1)0}^2 + b_{(t+1)1}^2 + b_{(t+1)2}^2 + \dots + b_{(t+1)t}^2) \\ &\geq \|a_t \mathbf{v} + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t\|_2^{2(k-1)} - O \left(|c_t|^{2(k-1)} \left(\frac{1}{n^\varepsilon |\beta|^{1/(k-2)}} + \frac{\log n}{\sqrt{n}} \right)^{2(k-1)} \right). \end{aligned} \quad (64)$$

For the first term on the righthand side of (64), using our induction hypothesis (60) and (61) that $|a_t| \lesssim |c_t|$, we have

$$\begin{aligned} \|a_t \mathbf{v} + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t\|_2^2 &= a_t^2 + b_t^2 + c_t^2 \|\boldsymbol{\xi}_t\|_2^2 + 2a_t b_t \langle \mathbf{v}, \mathbf{w}_t \rangle + 2a_t c_t \langle \mathbf{v}, \boldsymbol{\xi}_t \rangle + 2b_t c_t \langle \mathbf{w}_t, \boldsymbol{\xi}_t \rangle \\ &= \left(1 + O \left(\frac{\log n}{\sqrt{n}} + \frac{1}{n^{2\varepsilon} \beta^{2/(k-2)}} \right) \right) c_t^2. \end{aligned} \quad (65)$$

We get the following lower for c_{t+1} by plugging (65) into (64), and rearranging

$$|c_{t+1}| \geq \left(1 + O\left(\frac{\log n}{\sqrt{n}} + \frac{1}{n^{2\varepsilon}\beta^{2/(k-2)}}\right)\right) |c_t|^{k-1} \quad (66)$$

The claim that $|b_{(t+1)0}|, |b_{(t+1)1}|, \dots, |b_{(t+1)t}| \lesssim \max\{|c_{t+1}|/|\beta|^{(k-1)/(k-2)}, (\log n)^{k-1}|c_{t+1}|/n^{(k-1)/2}\}$ follows from combining (63) and (66). The claim that $|c_{t+1}| \geq n^\varepsilon \beta^{1/(k-2)} |a_{t+1}|$ follows from combining (62) and (66).

For (61), since τ_t is orthogonal to $\mathbf{v}^{\otimes(k-1)}, \tau_0, \tau_1, \dots, \tau_{t-1}$, Lemma 12 implies that conditioning on $\boldsymbol{\xi} = \mathbf{Z}[\mathbf{v}^{\otimes(k-1)}]$ and $\boldsymbol{\xi}_{s+1} = \mathbf{Z}[\tau_s]$ for $0 \leq s \leq t-1$, $\boldsymbol{\xi}_{t+1} = \mathbf{Z}[\tau_t]$ is an independent Gaussian vector, with each entry $\mathcal{N}(0, 1/n)$. By the standard concentration inequality, it holds that with probability $1 - e^{-c(\log n)^2}$, $\|\boldsymbol{\xi}_{t+1}\|_2 = 1 + O(\log n/\sqrt{n})$, and the projection of $\boldsymbol{\xi}_{t+1}$ on the span of $\{\mathbf{v}, \boldsymbol{\xi}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_t\}$ is bounded by $\log n/\sqrt{n}$. This finishes the proof of the induction (61).

Next, using (36) and (37) as input, we prove that for

$$t \geq 1 + \frac{1}{\varepsilon} \left(\frac{1}{2} - \frac{\log |\beta|}{(k-2) \log n} \right), \quad (67)$$

we have

$$\mathbf{y}_t = a_t \mathbf{v} + b_{t0} \boldsymbol{\xi} + b_{t1} \boldsymbol{\xi}_1 + \dots + b_{t(t-1)} \boldsymbol{\xi}_{t-1} + c_t \boldsymbol{\xi}_t, \quad (68)$$

such that

$$|a_t|, |b_{t0}|, |b_{t1}|, \dots, |b_{t(t-1)}| \lesssim |c_t| |\beta| \left(\frac{\log n}{\sqrt{n}} \right)^{k-1}. \quad (69)$$

Let $x_t = |a_t/c_t|$, then (60) implies that $x_t \leq 1/(n^\varepsilon |\beta|^{1/(k-2)})$. By taking the ratio of (62) and (66), we get

$$x_{t+1} \lesssim |\beta| \left(\frac{\log n}{\sqrt{n}} + x_t \right)^{k-1}. \quad (70)$$

there are two cases,

1. if $\log n/\sqrt{n} \lesssim x_t \leq 1/(n^\varepsilon |\beta|^{1/(k-2)})$, then

$$x_{t+1} \lesssim |\beta| x_t^{k-1} = x_t (|\beta|^{1/(k-2)} x_t)^{k-2} \leq x_t/n^\varepsilon;$$

2. If $x_t \lesssim \log n/\sqrt{n}$, then $|x_{t+1}| \lesssim |\beta| (\log n/\sqrt{n})^{k-1}$.

Since $x_1 = |a_1/c_1| \lesssim 1/(n^\varepsilon |\beta|^{1/(k-2)})$, we conclude that

$$x_t = |a_t/c_t| \lesssim |\beta| (\log n/\sqrt{n})^{k-1}, \quad \text{when } t \geq \frac{1}{\varepsilon} \left(\frac{1}{2} - \frac{\log |\beta|}{(k-2) \log n} \right). \quad (71)$$

In this regime, (63) implies that

$$\begin{aligned} & |b_{(t+1)0}|, |b_{(t+1)1}|, |b_{(t+1)2}|, \dots, |b_{(t+1)t}| \lesssim |\beta| |c_t|^{k-1} \left(\frac{|a_t|}{|c_t|} + \frac{\log n}{\sqrt{n}} \right)^{k-1} \\ & \lesssim |\beta| |c_t|^{k-1} \left(\frac{\log n}{\sqrt{n}} \right)^{k-1} \lesssim |c_{t+1}| |\beta| \left(\frac{\log n}{\sqrt{n}} \right)^{k-1}, \end{aligned}$$

where we used (66) in the last inequality. This finishes the proof of (69). Using (69), we can compute \mathbf{u}_t ,

$$\mathbf{u}_t = \frac{\mathbf{y}_t}{\|\mathbf{y}_t\|} = \frac{\boldsymbol{\xi}_t}{\|\boldsymbol{\xi}_t\|_2} + \mathcal{O}_{\mathbb{P}} \left(|\beta| \left(\frac{\log n}{\sqrt{n}} \right)^{k-1} \right),$$

where the error term is a vector of length bounded by $|\beta|(\log n/\sqrt{n})^{k-1}$. This finishes the proof of Theorem 1. ■

A.2 Proof of Corollarys 3 and 4

Proof [Proof of Corollary 3] According to the definition of $\boldsymbol{\xi}$ in (4) of Theorem 1, i.e. $\boldsymbol{\xi} = \mathbf{Z}[\mathbf{v}^{\otimes(k-1)}]$, is an n -dim vector, with each entry i.i.d. $\mathcal{N}(0, 1/n)$ Gaussian random variable. We see that

$$\langle \boldsymbol{\xi}, \mathbf{v} \rangle \stackrel{d}{=} \mathcal{N}(0, 1/n).$$

Especially with high probability we will have that $|\langle \boldsymbol{\xi}, \mathbf{v} \rangle| \lesssim \log n/\sqrt{n}$. Then we conclude from (5), with high probability it holds

$$\widehat{\beta} = \beta + \mathcal{O} \left(\frac{1}{\beta} + \frac{\log n}{\sqrt{n}} \right). \quad (72)$$

With the bound (112), we can replace $\langle \mathbf{a}, \mathbf{v} \rangle / (2\beta^2)$ on the righthand side of (4) by $\langle \mathbf{a}, \mathbf{v} \rangle / (2\widehat{\beta}^2)$, which gives an error

$$\left| \frac{\langle \mathbf{a}, \mathbf{v} \rangle}{2\beta^2} - \frac{\langle \mathbf{a}, \mathbf{v} \rangle}{2\widehat{\beta}^2} \right| = \mathcal{O} \left(|\langle \mathbf{a}, \mathbf{v} \rangle| \left(\frac{1}{|\beta|^4} + \frac{\log n}{|\beta|^3 \sqrt{n}} \right) \right).$$

Combining the above discussion together, we can rewrite (4) as

$$\langle \mathbf{a}, \widehat{\mathbf{v}} \rangle - \left(1 - \frac{1}{2\widehat{\beta}^2} \right) \langle \mathbf{a}, \mathbf{v} \rangle = \frac{\langle \mathbf{a}, \boldsymbol{\xi} \rangle - \langle \mathbf{a}, \mathbf{v} \rangle \langle \mathbf{v}, \boldsymbol{\xi} \rangle}{\beta} + \mathcal{O} \left(\frac{\log n}{\beta^2 \sqrt{n}} + \frac{(\log n)^{3/2}}{\beta^{3/2} n^{3/4}} + \frac{|\langle \mathbf{a}, \mathbf{v} \rangle|}{\beta^4} \right) \quad (73)$$

with high probability.

Again thanks to the definition of $\boldsymbol{\xi}$ in (4) of Theorem 1, i.e. $\boldsymbol{\xi} = \mathbf{Z}[\mathbf{v}^{\otimes(k-1)}]$, is an n -dim vector, with each entry i.i.d. $\mathcal{N}(0, 1/n)$ Gaussian random variable, we see that

$$\langle \mathbf{a}, \boldsymbol{\xi} \rangle - \langle \mathbf{a}, \mathbf{v} \rangle \langle \mathbf{v}, \boldsymbol{\xi} \rangle = \langle \mathbf{a} - \langle \mathbf{a}, \mathbf{v} \rangle \mathbf{v}, \boldsymbol{\xi} \rangle,$$

is a Gaussian random variable, with mean zero and variance

$$\mathbb{E}[\langle \mathbf{a} - \langle \mathbf{a}, \mathbf{v} \rangle \mathbf{v}, \boldsymbol{\xi} \rangle^2] = \frac{1}{n} \|\mathbf{a} - \langle \mathbf{a}, \mathbf{v} \rangle \mathbf{v}\|_2^2 = \frac{1}{n} \langle \mathbf{a}, (\mathbf{I}_n - \mathbf{v} \mathbf{v}^\top) \mathbf{a} \rangle = \frac{1 + o(1)}{n} \langle \mathbf{a}, (\mathbf{I}_n - \widehat{\mathbf{v}} \widehat{\mathbf{v}}^\top) \mathbf{a} \rangle.$$

This together with (112), (113) as well as our assumption (6)

$$\frac{\sqrt{n} \widehat{\beta}}{\sqrt{\langle \mathbf{a}, (\mathbf{I}_n - \widehat{\mathbf{v}} \widehat{\mathbf{v}}^\top) \mathbf{a} \rangle}} \left[\left(1 - \frac{1}{2\widehat{\beta}^2}\right)^{-1} \langle \mathbf{a}, \widehat{\mathbf{v}} \rangle - \langle \mathbf{a}, \mathbf{v} \rangle \right] \xrightarrow{d} \mathcal{N}(0, 1). \quad (74)$$

Under the same assumption, we have similar results for Cases 2, 3, 4, by simply changing (β, \mathbf{v}) in the righthand side of (4) and (5) to the corresponding expression. \blacksquare

Proof [Proof of Corollary 4] Given the significance level α , the asymptotic confidence intervals in Corollary 4 can be calculated from Corollary 3 by bounding the absolute values of the left hand sides of (7) at z_α . \blacksquare

A.3 Proof of Theorem 7

Proof [Proof of Theorem 7] We define an auxiliary iteration, $\mathbf{y}_0 = \mathbf{u}$ and

$$\mathbf{y}_{t+1} = \mathbf{X}[\mathbf{y}_t^{\otimes(k-1)}]. \quad (75)$$

Then we have that $\mathbf{u}_t = \mathbf{y}_t / \|\mathbf{y}_t\|_2$.

For index $\mathbf{j} = (j_1, j_2, \dots, j_{k-1}) \in \llbracket 1, r \rrbracket^{k-1}$. Let $\boldsymbol{\xi}_{\mathbf{j}} = \mathbf{Z}[\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \dots \otimes \mathbf{v}_{j_{k-1}}]$. Its entries

$$\boldsymbol{\xi}_{\mathbf{j}}(i) = \sum_{i_1, i_2, \dots, i_{k-1} \in \llbracket 1, n \rrbracket} \mathbf{Z}_{ii_1 i_2 \dots i_{k-1}} \mathbf{v}_{j_1}(i_1) \mathbf{v}_{j_2}(i_2) \dots \mathbf{v}_{j_{k-1}}(i_{k-1}),$$

are linear combination of Gaussian random variables, which is also Gaussian. These entries are i.i.d. Gaussian variables with mean zero and variance $1/n$,

$$\mathbb{E}[\boldsymbol{\xi}_{\mathbf{j}}(i)^2] = \sum_{i_1, i_2, \dots, i_{k-1} \in \llbracket 1, n \rrbracket} \mathbb{E}[\mathbf{Z}_{ii_1 i_2 \dots i_{k-1}}^2] \mathbf{v}_{j_1}(i_1)^2 \mathbf{v}_{j_2}(i_2)^2 \dots \mathbf{v}_{j_{k-1}}(i_{k-1})^2 = \frac{1}{n}.$$

We can compute \mathbf{y}_t iteratively:

$$\mathbf{y}_1 = \mathbf{X}[\mathbf{y}_0^{\otimes(k-1)}] = \sum_{j=1}^r \beta_j \langle \mathbf{y}_0, \mathbf{v}_j \rangle^{k-1} \mathbf{v}_j + \mathbf{Z}[\mathbf{y}_0^{\otimes(k-1)}]. \quad (76)$$

For the last term on the righthand side of (76), we can decompose $\mathbf{y}_0^{\otimes(k-1)}$ as a projection on $\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \dots \otimes \mathbf{v}_{j_{k-1}}$ for $\mathbf{j} \in \llbracket 1, r \rrbracket^{k-1}$, and its orthogonal part:

$$\mathbf{y}_0^{\otimes(k-1)} = \sum_{\mathbf{j}} \prod_{s=1}^{k-1} \langle \mathbf{y}_0, \mathbf{v}_{j_s} \rangle \mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \dots \otimes \mathbf{v}_{j_{k-1}} + \sqrt{1 - \left(\sum_{\mathbf{j}} \langle \mathbf{y}_0, \mathbf{v}_j \rangle^2 \right)^{(k-1)}} \boldsymbol{\tau}_0,$$

where the sum is over $\mathbf{j} \in \llbracket 1, r \rrbracket^{k-1}$, $\boldsymbol{\tau}_0 \in \otimes^k \mathbb{R}^n$ and $\|\boldsymbol{\tau}_0\|_2 = 1$. Let $\boldsymbol{\xi}_1 = \mathbf{Z}[\boldsymbol{\tau}_0]$. By our construction $\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}$ for any $\mathbf{j} \in \llbracket 1, r \rrbracket^{k-1}$ and $\boldsymbol{\tau}_0$ are othogonal to each other. Thanks to Lemma 12, conditioning on $\boldsymbol{\xi}_{\mathbf{j}} := \mathbf{Z}[\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}]$ for index $\mathbf{j} = (j_1, j_2, \dots, j_{k-1}) \in \llbracket 1, r \rrbracket^{k-1}$, $\boldsymbol{\xi}_1 = \mathbf{Z}[\boldsymbol{\tau}_0]$ has the same law as $\tilde{\mathbf{Z}}[\boldsymbol{\tau}_0]$, where $\tilde{\mathbf{Z}}$ is an independent copy of \mathbf{Z} . Since $\langle \boldsymbol{\tau}_0, \boldsymbol{\tau}_0 \rangle = 1$, $\boldsymbol{\xi}_1$ is a Gaussian vector with each entry $\mathcal{N}(0, 1/n)$. With those notations we can rewrite \mathbf{y}_1 as

$$\mathbf{y}_1 = \sum_{j=1}^r \beta_j \langle \mathbf{y}_0, \mathbf{v}_j \rangle^{k-1} \mathbf{v}_j + \sum_{\mathbf{j}} \prod_{s=1}^{k-1} \langle \mathbf{y}_0, \mathbf{v}_{j_s} \rangle \boldsymbol{\xi}_{\mathbf{j}} + \sqrt{1 - \left(\sum_{j=1}^r \langle \mathbf{y}_0, \mathbf{v}_j \rangle^2 \right)^{(k-1)}} \boldsymbol{\xi}_1. \quad (77)$$

In the following we show that:

Claim 15 *We can compute $\mathbf{y}_2, \mathbf{y}_3, \dots, \mathbf{y}_t$ inductively. The Gram-Schmidt orthonormalization procedure gives an orthogonal base of $\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}$ for $\mathbf{j} \in \llbracket 1, r \rrbracket^{k-1}$ and $\mathbf{y}_0^{\otimes(k-1)}, \mathbf{y}_1^{\otimes(k-1)}, \dots, \mathbf{y}_{t-1}^{\otimes(k-1)}$ as:*

$$\{\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}\}_{\mathbf{j} \in \llbracket 1, r \rrbracket^{k-1}}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}. \quad (78)$$

Let $\boldsymbol{\xi}_{\mathbf{j}} = \mathbf{Z}[\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}]$ for $\mathbf{j} = (j_1, j_2, \dots, j_{k-1}) \in \llbracket 1, r \rrbracket^{k-1}$, and $\boldsymbol{\xi}_{s+1} = \mathbf{Z}[\boldsymbol{\tau}_s]$ for $0 \leq s \leq t-1$. Conditioning on $\boldsymbol{\xi}_{\mathbf{j}} = \mathbf{Z}[\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}]$ for $\mathbf{j} = (j_1, j_2, \dots, j_{k-1}) \in \llbracket 1, r \rrbracket^{k-1}$ and $\boldsymbol{\xi}_{s+1} = \mathbf{Z}[\boldsymbol{\tau}_s]$ for $0 \leq s \leq t-2$, $\boldsymbol{\xi}_t = \mathbf{Z}[\boldsymbol{\tau}_{t-1}]$ is an independent Gaussian vector, with each entry $\mathcal{N}(0, 1/n)$. Then \mathbf{y}_t is in the following form

$$\mathbf{y}_t = a_t \mathbf{v}_t + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t, \quad (79)$$

where

$$a_t \mathbf{v}_t = a_{t1} \mathbf{v}_1 + a_{t2} \mathbf{v}_2 + \cdots + a_{tr} \mathbf{v}_r, \quad b_t \mathbf{w}_t = \sum_{\mathbf{j}} b_{t\mathbf{j}} \boldsymbol{\xi}_{\mathbf{j}} + b_{t1} \boldsymbol{\xi}_1 + \cdots + b_{t,t-1} \boldsymbol{\xi}_{t-1}, \quad (80)$$

and $\|\mathbf{v}_1\|_2, \|\mathbf{v}_2\|_2, \dots, \|\mathbf{v}_r\|_2, \|\mathbf{w}_t\|_2 = 1$.

Proof [Proof of Claim 15] The Claim 15 for $t = 1$ follows from (77). In the following, assuming Claim 15 holds for t , we prove it for $t+1$.

Conditioning on $\boldsymbol{\xi}_{\mathbf{j}} = \mathbf{Z}[\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}]$ for index $\mathbf{j} = (j_1, j_2, \dots, j_{k-1}) \in \llbracket 1, r \rrbracket^{k-1}$ and $\mathbf{Z}[\boldsymbol{\tau}_s] = \boldsymbol{\xi}_{s+1}$ for $0 \leq s \leq t-2$, Lemma 12 implies that $\boldsymbol{\xi}_t = \mathbf{Z}[\boldsymbol{\tau}_{t-1}]$ has the same law as $\tilde{\mathbf{Z}}[\boldsymbol{\tau}_{t-1}]$, where $\tilde{\mathbf{Z}}$ is an independent copy of \mathbf{Z} . Since $\boldsymbol{\tau}_{t-1}$ is orthogonal to $\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}$ for index $\mathbf{j} = (j_1, j_2, \dots, j_{k-1}) \in \llbracket 1, r \rrbracket^{k-1}$ and $\mathbf{Z}[\boldsymbol{\tau}_s] = \boldsymbol{\xi}_{s+1}$ for $0 \leq s \leq t-2$, $\boldsymbol{\xi}_t$ is an independent Gaussian random vector with each entry $\mathcal{N}(0, 1/n)$.

Let $\{\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}\}_{\mathbf{j} \in \llbracket 1, r \rrbracket^{k-1}}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_t$ be an orthogonal base for $\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}$ for $\mathbf{j} \in \llbracket 1, r \rrbracket^{k-1}$ and $\mathbf{y}_0^{\otimes(k-1)}, \mathbf{y}_1^{\otimes(k-1)}, \dots, \mathbf{y}_t^{\otimes(k-1)}$, obtained by the Gram-Schmidt orthonormalization procedure. More precisely, given those tensors $\{\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}\}_{\mathbf{j} \in \llbracket 1, r \rrbracket^{k-1}}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}$, we denote

$$\begin{aligned} b_{(t+1)\mathbf{j}} &= \langle \mathbf{y}_t^{\otimes(k-1)}, \mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}} \rangle, \quad \mathbf{j} = (j_1, j_2, \dots, j_{k-1}) \in \llbracket 1, r \rrbracket^{k-1}, \\ b_{(t+1)s} &= \langle \mathbf{y}_s^{\otimes(k-1)}, \boldsymbol{\tau}_{s-1} \rangle, \quad 1 \leq s \leq t, \quad c_{t+1} = \langle \mathbf{y}_t^{\otimes(k-1)}, \boldsymbol{\tau}_t \rangle \end{aligned} \quad (81)$$

and $\sum_j b_{(t+1)j} \mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}} + b_{(t+1)1} \boldsymbol{\tau}_0 + b_{(t+1)2} \boldsymbol{\tau}_1 + \cdots + b_{(t+1)t} \boldsymbol{\tau}_{t-1}$ is the projection of $\mathbf{y}_t^{\otimes(k-1)}$ on the span of $\{\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}\}_{j \in [1,r]^{k-1}}, \mathbf{y}_0^{\otimes(k-1)}, \mathbf{y}_1^{\otimes(k-1)}, \dots, \mathbf{y}_{t-1}^{\otimes(k-1)}$. Then we can write $\mathbf{y}_t^{\otimes(k-1)}$ in terms of the base (78)

$$\mathbf{y}_t^{\otimes(k-1)} = \sum_j b_{(t+1)j} \mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}} + b_{(t+1)1} \boldsymbol{\tau}_0 + b_{(t+1)2} \boldsymbol{\tau}_1 + \cdots + b_{(t+1)t} \boldsymbol{\tau}_{t-1} + c_{t+1} \boldsymbol{\tau}_t.$$

The recursion (75) implies that

$$\mathbf{y}_{t+1} = \sum_{j=1}^r \beta_j (a_{tj} + b_t \langle \mathbf{w}_t, \mathbf{v}_j \rangle + c_t \langle \boldsymbol{\xi}_t, \mathbf{v}_j \rangle)^{k-1} \mathbf{v}_j + b_{t+1} \mathbf{w}_{t+1} + c_{t+1} \mathbf{Z}[\boldsymbol{\tau}_t] \quad (82)$$

where

$$\begin{aligned} b_{t+1} \mathbf{w}_{t+1} &= \mathbf{Z} \left[\sum_j b_{(t+1)j} \mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}} + b_{(t+1)1} \boldsymbol{\tau}_0 + b_{(t+1)2} \boldsymbol{\tau}_1 + \cdots + b_{(t+1)t} \boldsymbol{\tau}_{t-1} \right] \\ &= \sum_j b_{(t+1)j} \boldsymbol{\xi}_j + b_{(t+1)1} \boldsymbol{\xi}_1 + b_{(t+1)2} \boldsymbol{\xi}_2 + \cdots + b_{(t+1)t} \boldsymbol{\xi}_t. \end{aligned} \quad (83)$$

Since $\boldsymbol{\tau}_t$ is orthogonal to $\{\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}\}_{j \in [1,r]^{k-1}}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}$, Lemma 12 implies that conditioning on $\boldsymbol{\xi}_j = \mathbf{Z}[\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}]$ for $\mathbf{j} = (j_1, j_2, \dots, j_{k-1}) \in [1, r]^{k-1}$ and $\boldsymbol{\xi}_{s+1} = \mathbf{Z}[\boldsymbol{\tau}_s]$ for $0 \leq s \leq t-1$, $\boldsymbol{\xi}_{t+1} = \mathbf{Z}[\boldsymbol{\tau}_t]$ is an independent Gaussian vector, with each entry $\mathcal{N}(0, 1/n)$. The above discussion gives us that

$$\mathbf{y}_{t+1} = a_{t+1} \mathbf{v}_{t+1} + b_{t+1} \mathbf{w}_{t+1} + c_{t+1} \boldsymbol{\xi}_{t+1}, \quad a_{t+1} = \sqrt{a_{(t+1)1}^2 + a_{(t+1)2}^2 + \cdots + a_{(t+1)r}^2}$$

and

$$a_{(t+1)j} = \beta_j (a_{tj} + b_t \langle \mathbf{w}_t, \mathbf{v}_j \rangle + c_t \langle \boldsymbol{\xi}_t, \mathbf{v}_j \rangle)^{k-1}, \quad 1 \leq j \leq r. \quad (84)$$

■

We recall that by our Assumption 5, that

$$1/\kappa \leq \left| \frac{\langle \mathbf{u}, \mathbf{v}_i \rangle}{\langle \mathbf{u}, \mathbf{v}_j \rangle} \right| \leq \kappa,$$

for all $1 \leq i, j \leq r$. If $j_* = \operatorname{argmax}_j \beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle^{k-2}$, it is necessary that $\beta_{j_*} \gtrsim \beta_1$, where the implicit constant depends on κ .

In the following, we study the case that $\langle \mathbf{u}, \mathbf{v}_{j_*} \rangle > 0$. The case $\langle \mathbf{u}, \mathbf{v}_{j_*} \rangle < 0$ can be proven in exactly the same way, by simply changing $(\beta, \mathbf{v}_{j_*})$ with $((-1)^k \beta, -\mathbf{v}_{j_*})$. We prove by induction

Claim 16 For any fixed time t , with probability at least $1 - O(e^{-c(\log n)^2})$ the following holds: for any $s \leq t$,

$$\begin{aligned} |a_{sj_*}| &\geq |a_{sj}|, \quad |a_s| \gtrsim |\beta_1| \left(\sum_j |b_{sj}| + |b_{s1}| + \cdots + |b_{s(s-1)}| \right), \\ |a_s| &\gtrsim n^\varepsilon \max\{\mathbf{1}(k \geq 3) |c_s / \beta_1^{1/(k-2)}|, |c_s / \sqrt{n}|\}, \end{aligned} \quad (85)$$

and for $\mathbf{j} = (j_1, j_2, \dots, j_{k-1}) \in \llbracket 1, r \rrbracket^{k-1}$

$$\begin{aligned} \|\boldsymbol{\xi}_j\|, \|\boldsymbol{\xi}_s\|_2 &= 1 + O(\log n / \sqrt{n}), \quad |\langle \mathbf{v}_j, \boldsymbol{\xi}_j \rangle|, |\langle \mathbf{a}, \boldsymbol{\xi}_j \rangle|, |\langle \mathbf{a}, \boldsymbol{\xi}_s \rangle| \lesssim \log n / \sqrt{n}. \\ \|\text{Proj}_{\text{Span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r, \{\boldsymbol{\xi}_j\}_{j \in \llbracket 1, r \rrbracket^{k-1}}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{s-1}\}}(\boldsymbol{\xi}_s)\|_2 &\lesssim \log n / \sqrt{n}. \end{aligned} \quad (86)$$

Proof [Proof of Claim 16] From (77), we have

$$\begin{aligned} \mathbf{y}_1 &= \sum_{j=1}^r \beta_j \langle \mathbf{y}_0, \mathbf{v}_j \rangle^{k-1} \mathbf{v}_j + \sum_j \prod_{s=1}^{k-1} \langle \mathbf{y}_0, \mathbf{v}_{j_s} \rangle \boldsymbol{\xi}_j + \sqrt{1 - \left(\sum_{j=1}^r \langle \mathbf{y}_0, \mathbf{v}_j \rangle^2 \right)^{(k-1)}} \boldsymbol{\xi}_1 \\ &= \sum_{j=1}^r a_{1j} \mathbf{v}_j + \sum_j b_{1j} \boldsymbol{\xi}_j + c_1 \boldsymbol{\xi}_1, \end{aligned}$$

where $a_{1j} = \beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle^{k-1}$ for $1 \leq j \leq r$, $b_{1j} = \prod_{s=1}^{k-1} \langle \mathbf{u}, \mathbf{v}_{j_s} \rangle$ for any index $\mathbf{j} = (j_1, j_2, \dots, j_{k-1})$ and $c_1 = \sqrt{1 - \left(\sum_{j=1}^r \langle \mathbf{u}, \mathbf{v}_j \rangle^2 \right)^{(k-1)}}$. Since $\boldsymbol{\xi}_j$ are independent Gaussian vectors with each entry mean zero and variance $1/n$, the concentration for chi-square distribution implies that $\|\boldsymbol{\xi}_j\|_2 = 1 + O(\log n / \sqrt{n})$ with probability $1 - e^{-c(\log n)^2}$. Since $j_* = \text{argmax}_j |\beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle^{k-2}|$, combining with our Assumption 5, it gives that $|a_{1j_*}| \geq |a_{1j}| / \kappa$. As a consequence, we also have that $|a_1| = \sqrt{a_{11}^2 + a_{12}^2 + \cdots + a_{1r}^2} \asymp |a_{1j_*}|$. Again using our Assumption 5

$$\sum_j |b_{1j}| \lesssim \left(\sum_{j=1}^r |\langle \mathbf{u}, \mathbf{v}_j \rangle| \right)^{k-1} \lesssim \sum_{j=1}^r |\langle \mathbf{u}, \mathbf{v}_j \rangle|^{k-1} \lesssim |a_{1j_*}| / \beta_{j_*} \lesssim |a_1| / \beta_1.$$

We can check that $|\beta_1^{1/(k-2)} a_1| \asymp |\beta_{j_*}^{1/(k-2)} a_{1j_*}| = |\beta_{j_*} \langle \mathbf{u}, \mathbf{v}_{j_*} \rangle^{k-2}|^{(k-1)/(k-2)} \geq n^\varepsilon \geq n^\varepsilon |c_1|$, and $|\sqrt{n} a_1| \asymp |\sqrt{n} a_{1j_*}| = |\beta_{j_*} \langle \mathbf{u}, \mathbf{v}_{j_*} \rangle^{k-2}| |\sqrt{n} \langle \mathbf{u}, \mathbf{v}_{j_*} \rangle| \gtrsim n^\varepsilon \geq n^\varepsilon |c_1|$. Moreover, conditioning on $\boldsymbol{\xi}_j = \mathbf{Z}[\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}]$ for $\mathbf{j} = (j_1, j_2, \dots, j_{k-1}) \in \llbracket 1, r \rrbracket^{k-1}$, Lemma 12 implies that $\boldsymbol{\xi}_1 = \mathbf{Z}[\boldsymbol{\tau}_0]$ is an independent Gaussian random vector with each entry $\mathcal{N}(0, 1/n)$. By the standard concentration inequality, it holds that with probability $1 - e^{-c(\log n)^2}$, $\|\boldsymbol{\xi}_1\|_2 = 1 + O(\log n / \sqrt{n})$, $|\langle \mathbf{a}, \boldsymbol{\xi}_1 \rangle|$ and the projection of $\boldsymbol{\xi}_1$ on the span of $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r, \{\boldsymbol{\xi}_j\}_{j \in \llbracket 1, r \rrbracket^{k-1}}\}$ is bounded by $\log n / \sqrt{n}$. So far we have proved that (36) and (37) hold for $t = 1$.

In the following, we assume that (85) and (86) hold for t , and prove it for $t + 1$. We recall from (80) and (84) that

$$a_{(t+1)j} = \beta_j (a_{tj} + b_t \langle \mathbf{w}_t, \mathbf{v}_j \rangle + c_t \langle \boldsymbol{\xi}_t, \mathbf{v}_j \rangle)^{k-1}, \quad b_t \mathbf{w}_t = \sum_j b_{tj} \boldsymbol{\xi}_j + b_{t1} \boldsymbol{\xi}_1 + \cdots + b_{tt-1} \boldsymbol{\xi}_{t-1}. \quad (87)$$

By our induction hypothesis, we have

$$|b_t \langle \mathbf{w}_t, \mathbf{v}_j \rangle| \lesssim \sum_j |b_{tj} \langle \boldsymbol{\xi}_j, \mathbf{v}_j \rangle| + |b_{t1} \langle \boldsymbol{\xi}_1, \mathbf{v}_j \rangle| + \cdots + |b_{t(t-1)} \langle \boldsymbol{\xi}_{t-1}, \mathbf{v}_j \rangle| \lesssim (\log n / \sqrt{n}) |a_t| / |\beta_1|, \quad (88)$$

and

$$|c_t \langle \boldsymbol{\xi}_t, \mathbf{v}_j \rangle| \lesssim (\log n / \sqrt{n}) |c_t| \lesssim (\log n) |a_t| / n^\varepsilon. \quad (89)$$

It follows from plugging (88) and (89) into (87), we get

$$a_{(t+1)j} = \beta_j (a_{tj} + b_t \langle \mathbf{w}_t, \mathbf{v}_j \rangle + c_t \langle \boldsymbol{\xi}_t, \mathbf{v}_j \rangle)^{k-1} = \beta_j (a_{tj} + O(\log n |a_t| / n^\varepsilon))^k \lesssim \beta_j |a_t|^{k-1},$$

and especially

$$a_{(t+1)j_*} = \beta_{j_*} (a_{tj_*} + O(\log n |a_{tj_*}| / n^\varepsilon))^k = (1 + O(\log n / n^\varepsilon)) \beta_{j_*} a_{tj_*}^{k-1}.$$

Therefore, we conclude that

$$|a_{(t+1)j_*}| \asymp |\beta_{j_*} a_{tj_*}^{k-1}| \asymp |\beta a_t^{k-1}|, \quad (90)$$

and

$$|a_{(t+1)j}| \lesssim \beta_j |a_t|^{k-1} \lesssim \beta_{j_*} |a_{tj_*}|^{k-1} \lesssim |a_{(t+1)j_*}|.$$

We recall from (81), $\sum_j b_{(t+1)j} \mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}} + b_{(t+1)1} \boldsymbol{\tau}_0 + b_{(t+1)2} \boldsymbol{\tau}_1 + \cdots + b_{(t+1)t} \boldsymbol{\tau}_{t-1}$ is the projection of $\mathbf{y}_t^{\otimes(k-1)}$ on the span of $\{\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}\}_{j \in [1,r]^{k-1}}, \mathbf{y}_0^{\otimes(k-1)}, \mathbf{y}_1^{\otimes(k-1)}, \dots, \mathbf{y}_{t-1}^{\otimes(k-1)}$. We also recall that $\{\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}\}_{j \in [1,r]^{k-1}}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}$ are obtained from $\{\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}\}_{j \in [1,r]^{k-1}}, \mathbf{y}_0^{\otimes(k-1)}, \mathbf{y}_1^{\otimes(k-1)}, \dots, \mathbf{y}_{t-1}^{\otimes(k-1)}$ by the Gram-Schmidt orthonormalization procedure. So we have that the span of vectors $\{\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}\}_{j \in [1,r]^{k-1}}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}$ is the same as the span of the tensors $\{\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}\}_{j \in [1,r]^{k-1}}, \mathbf{y}_0^{\otimes(k-1)}, \mathbf{y}_1^{\otimes(k-1)}, \dots, \mathbf{y}_{t-1}^{\otimes(k-1)}$, which is contained in the span of the tensors $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r, \mathbf{w}_t, \mathbf{y}_0, \dots, \mathbf{y}_{t-1}\}^{\otimes(k-1)}$. Moreover from the relation (79) and (80), one can see that the span of $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r, \mathbf{w}_t, \mathbf{y}_0, \dots, \mathbf{y}_{t-1}\}$ is the same as the span of $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r, \{\boldsymbol{\xi}_j\}_{j \in [1,r]^{k-1}}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}\}$. It follows that

$$\begin{aligned} |b_{t+1}| &\lesssim \sqrt{\sum_j b_{(t+1)j}^2 + b_{(t+1)1}^2 + b_{(t+1)2}^2 + \cdots + b_{(t+1)t}^2} \\ &= \|\text{Proj}_{\text{Span}\{\{\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}\}_{j \in [1,r]^{k-1}}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}\}}(a_t \mathbf{v}_t + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t)^{\otimes(k-1)}\|_2 \\ &\leq \|\text{Proj}_{\text{Span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r, \mathbf{w}_t, \mathbf{y}_0, \dots, \mathbf{y}_{t-1}\}^{\otimes(k-1)}}(a_t \mathbf{v}_t + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t)^{\otimes(k-1)}\|_2 \\ &\leq \|\text{Proj}_{\text{Span}\{\mathbf{v}, \mathbf{w}_t, \mathbf{y}_0, \dots, \mathbf{y}_{t-1}\}}(a_t \mathbf{v}_t + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t)\|_2^{k-1} \\ &= \|a_t \mathbf{v}_t + b_t \mathbf{w}_t + c_t \text{Proj}_{\text{Span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r, \{\boldsymbol{\xi}_j\}_{j \in [1,r]^{k-1}}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}\}}(\boldsymbol{\xi}_t)\|_2^{k-1} \\ &\lesssim \left(|a_t| + |b_t| + \frac{\log n |c_t|}{\sqrt{n}} \right)^{k-1} \lesssim |a_t|^{k-1} \lesssim |a_{t+1}| / \beta_1, \end{aligned} \quad (91)$$

where in the first line we used (83), and in the last line of (91) we used our induction hypothesis that $\|\text{Proj}_{\text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r, \{\boldsymbol{\xi}_j\}_{j \in [1, r]^{k-1}}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}\}}(\boldsymbol{\xi}_t)\|_2 \lesssim \log n / \sqrt{n}$.

Finally we estimate c_{t+1} . We recall from (81), the coefficient c_{t+1} is the remainder of $\mathbf{y}_t^{\otimes(k-1)}$ after projecting on $\{\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \dots \otimes \mathbf{v}_{j_{k-1}}\}_{j \in [1, r]^{k-1}}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}$. It is bounded by the remainder of $\mathbf{y}_t^{\otimes(k-1)}$ after projecting on $\{\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \dots \otimes \mathbf{v}_{j_{k-1}}\}_{j \in [1, r]^{k-1}}$,

$$|c_{t+1}| \leq \|\mathbf{y}_t^{\otimes(k-1)} - a_t^{k-1} \mathbf{v}_t^{\otimes(k-1)}\|_2 = \|(a_t \mathbf{v}_t + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t)^{\otimes(k-1)} - a_t^{k-1} \mathbf{v}_t^{\otimes(k-1)}\|_2.$$

The difference $(a_t \mathbf{v}_t + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t)^{\otimes(k-1)} - a_t^{k-1} \mathbf{v}_t^{\otimes(k-1)}$ is a sum of terms in the following form,

$$\boldsymbol{\eta}_1 \otimes \boldsymbol{\eta}_2 \otimes \dots \otimes \boldsymbol{\eta}_{k-1}, \quad (92)$$

where $\boldsymbol{\eta}_1, \boldsymbol{\eta}_2, \dots, \boldsymbol{\eta}_{k-1} \in \{a_t \mathbf{v}_t, b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t\}$, and at least one of them is $b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t$. We notice that by our induction hypothesis, $\|b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t\|_2 \lesssim |b_t| \|\mathbf{w}_t\|_2 + |c_t| \|\boldsymbol{\xi}_t\|_2 \lesssim |b_t| + |c_t|$. For the L_2 norm of (92), each copy of $a_t \mathbf{v}_t$ contributes a_t and each copy of $b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t$ contributes a factor $|b_t| + |c_t|$. We conclude that

$$|c_{t+1}| \leq \|(a_t \mathbf{v}_t + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t)^{\otimes(k-1)} - a_t^{k-1} \mathbf{v}_t^{\otimes(k-1)}\|_2 \lesssim \sum_{r=1}^{k-1} |a_t|^{k-1-r} (|b_t| + |c_t|)^r. \quad (93)$$

Combining with (90) that $|a_{t+1}| \asymp |\beta_1| |a_t|^{k-1}$, we divide both sides of (93) by $|\beta_1| |a_t|^{k-1}$,

$$\frac{|c_{t+1}|}{|a_{t+1}|} \lesssim \frac{1}{|\beta_1|} \sum_{r=1}^{k-1} \left(\frac{|b_t|}{|a_t|} + \frac{|c_t|}{|a_t|} \right)^r \lesssim \frac{1}{|\beta_1|} \sum_{r=1}^{k-1} \left(\frac{1}{|\beta_1|} + \frac{|c_t|}{|a_t|} \right)^r \quad (94)$$

There are three cases:

1. If $|c_t|/|a_t| \geq 1$, then

$$\frac{|c_{t+1}|}{|a_{t+1}|} \lesssim \frac{1}{|\beta_1|} \sum_{r=1}^{k-1} \left(\frac{1}{|\beta_1|} + \frac{|c_t|}{|a_t|} \right)^r \lesssim \frac{1}{|\beta_1|} \left(\frac{|c_t|}{|a_t|} \right)^{k-1}.$$

If $k = 2$, then $|c_{t+1}|/|a_{t+1}| \lesssim (|c_t|/|a_t|)/n^\varepsilon$. If $k \geq 2$, by our induction hypothesis $|c_t|/|a_t| \lesssim \beta_1^{1/(k-2)}/n^\varepsilon$. Especially, $(|c_t|/|a_t|)^{k-2}/|\beta_1| \lesssim 1/n^\varepsilon$. We still get that $|c_{t+1}|/|a_{t+1}| \lesssim (|c_t|/|a_t|)/n^\varepsilon$.

2. If $1/|\beta_1| \lesssim |c_t|/|a_t| \leq 1$, then

$$\frac{|c_{t+1}|}{|a_{t+1}|} \lesssim \frac{1}{|\beta_1|} \sum_{r=1}^{k-1} \left(\frac{1}{|\beta_1|} + \frac{|c_t|}{|a_t|} \right)^r \lesssim \frac{1}{|\beta_1|} \left(\frac{|c_t|}{|a_t|} \right) \lesssim \frac{1}{n^\varepsilon} \left(\frac{|c_t|}{|a_t|} \right).$$

3. Finally for $|c_t|/|a_t| \lesssim 1/|\beta_1|$, we will have

$$\frac{|c_{t+1}|}{|a_{t+1}|} \lesssim \frac{1}{|\beta_1|} \sum_{r=1}^{k-1} \left(\frac{1}{|\beta_1|} + \frac{|c_t|}{|a_t|} \right)^r \lesssim \frac{1}{|\beta_1|} \left(\frac{1}{|\beta_1|} \right) \lesssim \frac{1}{|\beta_1|^2}.$$

In all these cases we have $|c_{t+1}|/|a_{t+1}| \lesssim \min\{\sqrt{n}, \mathbf{1}(k \geq 3)|\beta_1|^{1/(k-2)}\}/n^\varepsilon$. This finishes the proof of the induction (85).

For (86), since $\boldsymbol{\tau}_t$ is orthogonal to $\{\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}\}_{j \in \llbracket 1, r \rrbracket^{k-1}}, \boldsymbol{\tau}_0, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_{t-1}$, Lemma 12 implies that conditioning on $\boldsymbol{\xi}_j = \mathbf{Z}[\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \cdots \otimes \mathbf{v}_{j_{k-1}}]$ for index $\mathbf{j} = (j_1, j_2, \dots, j_{k-1}) \in \llbracket 1, r \rrbracket^{k-1}$ and $\boldsymbol{\xi}_{s+1} = \mathbf{Z}[\boldsymbol{\tau}_s]$ for $0 \leq s \leq t-1$, $\boldsymbol{\xi}_{t+1} = \mathbf{Z}[\boldsymbol{\tau}_t]$ is an independent Gaussian vector, with each entry $\mathcal{N}(0, 1/n)$. By the standard concentration inequality, it holds that with probability $1 - e^{c(\log n)^2}$, $\|\boldsymbol{\xi}_{t+1}\|_2 = 1 + O(\log n/\sqrt{n})$, $|\langle \mathbf{a}, \boldsymbol{\xi}_{t+1} \rangle|$ and the projection of $\boldsymbol{\xi}_{t+1}$ on the span of $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r, \{\boldsymbol{\xi}_j\}_{j \in \llbracket 1, r \rrbracket^{k-1}}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}\}$ is bounded by $\log n/\sqrt{n}$. This finishes the proof of the induction (86). \blacksquare

Next, using (85) and (86) as input, we prove that for

$$t \geq 1 + \frac{1}{\varepsilon} \left(\frac{1}{2} + \frac{2 \log |\beta_1|}{\log n} \right) + \frac{\log \log(\sqrt{n}|\beta_1|)}{\log(k-1)} \quad (95)$$

we have

$$\mathbf{y}_t = \sum_{j=1}^r a_{tj} \mathbf{v}_j + \sum_{\mathbf{j}} b_{tj} \boldsymbol{\xi}_{\mathbf{j}} + b_{t1} \boldsymbol{\xi}_1 + \cdots + b_{t,t-1} \boldsymbol{\xi}_{t-1} + c_t \boldsymbol{\xi}_t, \quad (96)$$

such that

$$\begin{aligned} |a_{tj}| &\lesssim \left(\frac{\log n}{\sqrt{n}} \frac{1}{|\beta_1|} \right)^{k-1} |a_{tj_*}|, \quad j \neq j_*, \\ b_{t(j_*, j_*, \dots, j_*)} &= \frac{a_{tj_*}}{\beta_{j_*}} + O\left(\frac{\log n |a_t|}{|\beta_1|^2 \sqrt{n}} \right), |b_{(t+1)j_*}| \lesssim \frac{\log n}{\sqrt{n} |\beta_1|^2} |a_{tj_*}|, \quad \mathbf{j}_* = (j_*, j_*, \dots, j_*), \\ |b_{t1}|, |b_{t2}|, \dots, |b_{t,t-1}| &\lesssim \frac{(\log n)^{1/2} |a_t|}{|\beta_1|^{3/2} n^{1/4}}, \quad |c_t| \lesssim |a_t|/\beta_1^2 \end{aligned} \quad (97)$$

Let $x_t = |c_t/a_t| \leq n^{-\varepsilon} |\beta|^{1/(k-2)}$, and $r_t = \max_{j \neq j_*} (\beta_j^{1/(k-2)} a_{tj}) / (\beta_{j_*}^{1/(k-2)} a_{tj_*})$. For $t = 1$, our Assumption 6 implies that

$$\begin{aligned} \beta_j^{1/(k-2)} a_{1j} &\leq (\beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle^{k-2})^{(k-1)/(k-2)} \\ &\leq ((1 - 1/\kappa) \beta_{j_*} \langle \mathbf{u}, \mathbf{v}_{j_*} \rangle^{k-2})^{(k-1)/(k-2)} \leq (1 - 1/\kappa) \beta_{j_*}^{1/(k-2)} a_{1j_*}. \end{aligned}$$

Thus we have that $r_1 \leq (1 - 1/\kappa)$. We recall from (87)

$$\begin{aligned} \beta_j^{1/(k-2)} a_{(t+1)j} &= \left(\beta_j^{1/k-2} (a_{tj} + b_t \langle \mathbf{w}_t, \mathbf{v}_j \rangle + c_t \langle \boldsymbol{\xi}_t, \mathbf{v}_j \rangle) \right)^{k-1} \\ &= \left(\beta_j^{1/k-2} \left(a_{tj} + O\left(|a_t| \frac{\log n (1/|\beta_1| + x_t)}{\sqrt{n}} \right) \right) \right)^{k-1}, \end{aligned}$$

where we used (85) and (86). Thus it follows that

$$\begin{aligned} r_{t+1} &= \max_{j \neq j_*} \left(\frac{\beta_j^{1/k-2} (a_{tj} + O(|a_t| \log n (1/|\beta_1| + x_t)/\sqrt{n}))}{\beta_{j_*}^{1/k-2} (a_{tj_*} + O(|a_t| \log n (1/|\beta_1| + x_t)/\sqrt{n}))} \right)^{k-1} \\ &\leq \left(\frac{r_t + O(\log n (1/|\beta_1| + x_t)/\sqrt{n})}{1 + O(\log n (1/|\beta_1| + x_t)/\sqrt{n})} \right)^{k-1} \end{aligned} \quad (98)$$

For x_t , (94) implies

$$x_{t+1} \lesssim \frac{1}{|\beta_1|} \sum_{r=1}^{k-1} \left(\frac{1}{|\beta_1|} + x_t \right)^r, \quad (99)$$

from the discussion after (94), we have that either $x_{t+1} \lesssim 1/|\beta_1|^2$, or $x_{t+1} \lesssim x_t/n^\varepsilon$. Since $x_1 = |c_1/a_1| \lesssim n^{1/2-\varepsilon}$, and $r_1 \leq (1 - 1/\kappa)$ we conclude from (98) and (99) that

$$x_t = |c_t/a_t| \lesssim 1/\beta_1^2, \quad r_t \lesssim (\log n / (|\beta_1| \sqrt{n}))^{k-1}, \quad (100)$$

when

$$t \geq \frac{1}{\varepsilon} \left(\frac{1}{2} + \frac{2 \log |\beta_1|}{\log n} \right) + \frac{\log \log(\sqrt{n} |\beta_1|)}{\log(k-1)}.$$

To derive the upper bound of $b_{t1}, b_{t2}, \dots, b_{t(t-1)}$, we use (91).

$$\begin{aligned} &\sum_j b_{(t+1)j}^2 + b_{(t+1)1}^2 + b_{(t+1)2}^2 + \dots + b_{(t+1)t}^2 \\ &\leq \|a_t \mathbf{v}_t + b_t \mathbf{w}_t + c_t \text{Proj}_{\text{Span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r, \{\boldsymbol{\xi}_j\}_{j \in [1, r]^{k-1}}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}\}}(\boldsymbol{\xi}_t)\|_2^{2(k-1)} \\ &= \left(a_t^2 + O \left(|a_t| (|b_t| + |c_t|) \frac{\log n}{\sqrt{n}} + \left(|b_t| + |c_t| \frac{\log n}{\sqrt{n}} \right)^2 \right) \right)^{k-1}, \end{aligned} \quad (101)$$

where we used (86). The first term $b_{(t+1)j}$ is the projection of $\mathbf{y}_t^{\otimes(k-1)}$ on $\mathbf{v}_{j_1} \otimes \mathbf{v}_{j_2} \otimes \dots \otimes \mathbf{v}_{j_{k-1}}$,

$$b_{(t+1)j} = \prod_{s=1}^{k-1} \langle a_t \mathbf{v}_t + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t, \mathbf{v}_{j_s} \rangle = \prod_{s=1}^{k-1} \left(a_{tj_s} + O \left(\frac{\log n (|b_t| + |c_t|)}{\sqrt{n}} \right) \right), \quad (102)$$

and

$$\begin{aligned} \sum_j b_{(t+1)j}^2 &= \left(\sum_{s=1}^{k-1} |\langle a_t \mathbf{v}_t + b_t \mathbf{w}_t + c_t \boldsymbol{\xi}_t, \mathbf{v}_{j_s} \rangle|^2 \right)^k \\ &= \left(a_t^2 + O \left(|a_t| (|b_t| + |c_t|) \frac{\log n}{\sqrt{n}} + \left(|b_t| + |c_t| \frac{\log n}{\sqrt{n}} \right)^2 \right) \right)^{k-1}, \end{aligned} \quad (103)$$

where we used (37) that $|\langle \boldsymbol{\xi}_j, \mathbf{v}_j \rangle|, |\langle \boldsymbol{\xi}_1, \mathbf{v}_j \rangle|, \dots, |\langle \boldsymbol{\xi}_t, \mathbf{v}_j \rangle| \lesssim \log n / \sqrt{n}$. Now we can take difference of (101) and (103), and use that $|b_t| \lesssim |a_t|/|\beta_1|$ from (85) and $|c_t| \lesssim |a_t|/|\beta_1|^2$ from (100),

$$b_{(t+1)1}^2 + b_{(t+1)2}^2 + \dots + b_{(t+1)t}^2 \lesssim a_t^{2(k-1)} \frac{\log n}{|\beta| \sqrt{n}}. \quad (104)$$

Using (102) and (100), we get that

$$\begin{aligned} b_{(t+1)\mathbf{j}_*} &= a_{t\mathbf{j}_*}^{k-1} \left(1 + \mathcal{O} \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right) \right), \quad \mathbf{j}_* = (j_*, j_*, \dots, j_*) \\ |b_{(t+1)\mathbf{j}}| &\lesssim \frac{\log n}{\sqrt{n}|\beta_1|} |a_{t\mathbf{j}_*}|^{k-1}, \quad \mathbf{j} \neq \mathbf{j}_*. \end{aligned} \quad (105)$$

From (87), (90) and (100), we have that

$$\begin{aligned} a_{(t+1)\mathbf{j}_*} &= \beta_{j_*} b_{(t+1)\mathbf{j}_*} = \beta_{j_*} a_{t\mathbf{j}_*}^{k-1} \left(1 + \mathcal{O} \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right) \right), \\ |a_{(t+1)\mathbf{j}}| &\lesssim \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right)^{k-1} |a_{(t+1)\mathbf{j}_*}|, \quad \mathbf{j} \neq \mathbf{j}_*. \end{aligned} \quad (106)$$

Using the above relation, we can simplify (104) and (105) as

$$|b_{(t+1)1}|, |b_{(t+1)2}|, \dots, |b_{(t+1)t}| \lesssim \frac{(\log n)^{1/2} |a_{t+1}|}{|\beta_1|^{3/2} n^{1/4}}.$$

and

$$\begin{aligned} b_{(t+1)\mathbf{j}_*} &= \frac{a_{(t+1)\mathbf{j}_*}}{\beta_{j_*}} \left(1 + \mathcal{O} \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right) \right), \\ |b_{(t+1)\mathbf{j}}| &\lesssim \frac{\log n}{\sqrt{n}|\beta_1|^2} |a_{(t+1)\mathbf{j}_*}|, \quad \mathbf{j} \neq \mathbf{j}_*. \end{aligned}$$

This finishes the proof of (97).

With the expression (97), we can process to prove our main results (11) and (12). Thanks to (86) and (96), for t satisfies (95), we have that with probability at least $1 - \mathcal{O}(e^{-c(\log n)^2})$

$$\|\mathbf{y}_t\|_2^2 = a_{t\mathbf{j}_*}^2 \left(1 + \frac{1}{\beta_{j_*}^2} + \frac{2\langle \mathbf{v}_{j_*}, \boldsymbol{\xi}_{j_*} \rangle}{\beta_{j_*}} + \mathcal{O} \left(\frac{\log n}{\sqrt{n}} \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right)^{k-1} + \frac{\log n}{\beta_1^2 \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta_1|^{3/2} n^{3/4}} + \frac{1}{\beta_1^4} \right) \right) \quad (107)$$

where $\mathbf{j}_* = (j_*, j_*, \dots, j_*)$. By rearranging it we get

$$1/\|\mathbf{y}_t\|_2 = a_{t\mathbf{j}_*}^2 \left(1 - \frac{1}{2\beta_{j_*}^2} - \frac{2\langle \mathbf{v}_{j_*}, \boldsymbol{\xi}_{j_*} \rangle}{\beta_{j_*}} + \mathcal{O} \left(\frac{\log n}{\sqrt{n}} \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right)^{k-1} + \frac{\log n}{\beta_1^2 \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta_1|^{3/2} n^{3/4}} + \frac{1}{\beta_1^4} \right) \right) \quad (108)$$

We can take the inner product $\langle \mathbf{a}, \mathbf{y}_t \rangle$, and multiply (108)

$$\begin{aligned} \langle \mathbf{a}, \mathbf{u}_t \rangle &= \frac{\langle \mathbf{a}, \mathbf{y}_t \rangle}{\|\mathbf{y}_t\|_2} = \text{sgn}(a_{tj_*}) \left(\left(1 - \frac{1}{2\beta_{j_*}^2} \right) \langle \mathbf{a}, \mathbf{v}_{j_*} \rangle + \frac{\langle \mathbf{a}, \boldsymbol{\xi}_{j_*} \rangle - \langle \mathbf{a}, \mathbf{v}_{j_*} \rangle \langle \mathbf{v}_{j_*}, \boldsymbol{\xi}_{j_*} \rangle}{\beta} \right) \\ &\quad + \text{O}_{\mathbb{P}} \left(\frac{\log n}{\sqrt{n}} \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right)^{k-1} + \frac{\log n}{\beta_1^2 \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta_1|^{3/2} n^{3/4}} + \frac{1}{\beta_1^4} \right), \end{aligned}$$

where we used (37) that with high probability $|\langle \mathbf{a}, \boldsymbol{\xi}_j \rangle|, |\langle \mathbf{a}, \boldsymbol{\xi}_s \rangle|$ for $1 \leq s \leq t$ are bounded by $\log n / \sqrt{n}$. This finishes the proof of (11). For $\hat{\beta}$ in (12), we have that

$$\mathbf{X}[\mathbf{u}_t^{\otimes k}] = \frac{\mathbf{X}[\mathbf{y}_t^{\otimes k}]}{\|\mathbf{y}_t\|_2^k} = \frac{\langle \mathbf{y}_t, \mathbf{X}[\mathbf{y}_t^{\otimes(k-1)}] \rangle}{\|\mathbf{y}_t\|_2^k} = \frac{\langle \mathbf{y}_t, \mathbf{y}_{t+1} \rangle}{\|\mathbf{y}_t\|_2^k}. \quad (109)$$

Thanks to (100), (106) and (86), for t satisfies (95), with probability at least $1 - \text{O}(e^{-c(\log n)^2})$, we can write the first term on the righthand side of (57), we have

$$\mathbf{y}_{t+1} = \sum_j a_{(t+1)j} \mathbf{v}_j + \sum_j b_{(t+1)j} \boldsymbol{\xi}_j + b_{(t+1)1} \boldsymbol{\xi}_1 + \cdots + b_{(t+1)t} \boldsymbol{\xi}_t + c_{t+1} \boldsymbol{\xi}_{t+1},$$

where $|c_{t+1}| \lesssim |a_t|^{k-1} / \beta^2$,

$$\begin{aligned} a_{(t+1)j_*} &= \beta_{j_*} b_{(t+1)j_*} = \beta_{j_*} a_{tj_*}^{k-1} \left(1 + \text{O} \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right) \right), \\ |a_{(t+1)j}| &\lesssim \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right)^{k-1} |a_{(t+1)j_*}|, \quad j \neq j_* \end{aligned}$$

and

$$\begin{aligned} b_{(t+1)j_*} &= \frac{a_{(t+1)j_*}}{\beta_{j_*}} \left(1 + \text{O} \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right) \right), \\ |b_{(t+1)j}| &\lesssim \frac{\log n}{\sqrt{n}|\beta_1|^2} |a_{(t+1)j_*}|, \quad j \neq j_*, \\ |b_{(t+1)1}|, |b_{(t+1)2}| + \cdots + |b_{(t+1)t}| &\lesssim a_t^{k-1} \frac{(\log n)^{1/2}}{|\beta|^{1/2} n^{1/4}}. \end{aligned}$$

From the discussion above, combining with (96) and (97) with straightforward computation, we have

$$\langle \mathbf{y}_t, \mathbf{y}_{t+1} \rangle = \beta_{j_*} a_{tj_*}^k \left(1 + \frac{1}{\beta_{j_*}^2} + \frac{(k+1)\langle \boldsymbol{\xi}_{j_*}, \mathbf{v}_{j_*} \rangle}{\beta_{j_*}} + \text{O} \left(\frac{\log n}{\sqrt{n}} \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right)^{k-1} + \frac{\log n}{\beta_1^2 \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta_1|^{3/2} n^{3/4}} \right) \right). \quad (110)$$

By plugging (108) and (110) into (109), we get

$$\begin{aligned} \mathbf{X}[\mathbf{u}_t^{\otimes k}] &= \text{sgn}(a_{tj_*}^k) \left(\beta_{j_*} + \langle \boldsymbol{\xi}_{j_*}, \mathbf{v}_{j_*} \rangle - \frac{k/2 - 1}{\beta_{j_*}} \right) \\ &\quad + \text{O} \left(\frac{\log n}{\sqrt{n}} \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right)^{k-1} + \frac{\log n}{|\beta_1| \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta_1|^{1/2} n^{3/4}} + \frac{1}{|\beta_1|^3} \right) \end{aligned}$$

Since by our assumption, in Case 1 we have that $\beta_{j^*} > 0$. Thanks to (106) $a_{t+1j^*} = \beta a_{tj^*}^{k-1}(1 + o(1))$, especially a_{t+1j^*} and a_{tj^*} are of the same sign. In the case $\langle \mathbf{u}, \mathbf{v}_{j^*} \rangle > 0$, we have $a_{1j^*} = \beta \langle \mathbf{u}, \mathbf{v}_{j^*} \rangle^{k-1} > 0$. We conclude that $a_{tj^*} > 0$. Therefore $\text{sgn}(\mathbf{X}[\mathbf{u}_t^{\otimes k}]) = \text{sgn}(a_{tj^*})^k = +$, and it follows that

$$\begin{aligned} \mathbf{X}[\mathbf{u}_t^{\otimes k}] &= \beta_{j^*} + \langle \boldsymbol{\xi}_{j^*}, \mathbf{v}_{j^*} \rangle - \frac{k/2 - 1}{\beta_{j^*}} \\ &+ \text{O} \left(\frac{\log n}{\sqrt{n}} \left(\frac{\log n}{\sqrt{n}|\beta_1|} \right)^{k-1} + \frac{\log n}{|\beta_1|\sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta_1|^{1/2}n^{3/4}} + \frac{1}{|\beta_1|^3} \right) \end{aligned}$$

This finishes the proof of (12). The Cases 2, 3, 4, by simply changing $(\beta_{j^*}, \mathbf{v}_{j^*})$ in the righthand side of (11) and (12) to the corresponding limit. ■

A.4 Proof of Theorem 9

Proof [Proof of Theorem 9] We first prove (15). If \mathbf{u} is uniformly distributed over the unit sphere, then it has the same law as $\boldsymbol{\eta}/\|\boldsymbol{\eta}\|_2$, where $\boldsymbol{\eta}$ is an n -dim standard Gaussian vector, with each entry $\mathcal{N}(0, 1)$. With this notation

$$|\beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle^{k-2}| = |\beta_j \langle \boldsymbol{\eta}, \mathbf{v}_j \rangle^{k-2}| / \|\boldsymbol{\eta}\|_2^{k-2}, \quad (111)$$

and we can rewrite $\mathbb{P}(i = \text{argmax}_j |\beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle^{k-2}|)$ as

$$\mathbb{P}(i = \text{argmax}_j |\beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle^{k-2}|) = \mathbb{P}(i = \text{argmax}_j |\beta_j \langle \boldsymbol{\eta}, \mathbf{v}_j \rangle^{k-2}|).$$

Since $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$ are orthonormal vectors, $\langle \mathbf{v}_1, \boldsymbol{\eta} \rangle, \langle \mathbf{v}_2, \boldsymbol{\eta} \rangle, \dots, \langle \mathbf{v}_r, \boldsymbol{\eta} \rangle$ are independent standard Gaussian random variables. Then we have

$$\begin{aligned} p_i &= \mathbb{P}(i = \text{argmax}_j |\beta_j \langle \boldsymbol{\eta}, \mathbf{v}_j \rangle^{k-2}|) = \mathbb{P}(|\beta_i/\beta_\ell|^{1/k-2} \langle \boldsymbol{\eta}, \mathbf{v}_i \rangle \geq |\langle \boldsymbol{\eta}, \mathbf{v}_\ell \rangle|, \text{ for all } i \neq \ell) \\ &= \int_0^\infty \sqrt{\frac{2}{\pi}} e^{-x^2/2} \left(\prod_{\ell \neq i} \int_0^{\left(\frac{|\beta_i|}{|\beta_\ell|}\right)^{\frac{1}{k-2}} x} \sqrt{\frac{2}{\pi}} e^{-y^2/2} dy \right) dx. \end{aligned}$$

This gives (15). Using the fact we can rewrite \mathbf{u} as $\boldsymbol{\eta}/\|\boldsymbol{\eta}\|_2$, we have that with probability $1 - \text{O}(1/\sqrt{\kappa})$,

$$1/\sqrt{\kappa n} \leq |\langle \mathbf{u}, \mathbf{v}_i \rangle| \leq \sqrt{\kappa/n},$$

for all $1 \leq i \leq r$. Thus Assumption (5) holds, and especially,

$$\max_j |\beta_j \langle \mathbf{u}, \mathbf{v}_j \rangle^{k-2}| \geq |\beta_1 \langle \mathbf{u}, \mathbf{v}_1 \rangle^{k-2}| \geq |\beta_1 (1/\sqrt{\kappa n})^{k-2}| \gtrsim n^\varepsilon.$$

Theorem 9 then follows directly from Theorem 7. ■

A.5 Proof of Corollaries 8, 10 and 11

Proof [Proof of Corollary 8] According to the definition of $\boldsymbol{\xi}$ in (11) of Theorem 7, i.e. $\boldsymbol{\xi} = \mathbf{Z}[\mathbf{v}_{j_*}^{\otimes(k-1)}]$, is an n -dim vector, with each entry i.i.d. $\mathcal{N}(0, 1/n)$ Gaussian random variable. We see that

$$\langle \boldsymbol{\xi}, \mathbf{v} \rangle \stackrel{d}{=} \mathcal{N}(0, 1/n).$$

Especially with high probability we will have that $|\langle \boldsymbol{\xi}, \mathbf{v} \rangle| \lesssim \log n / \sqrt{n}$. Then we conclude from (12), with high probability it holds

$$\widehat{\beta} = \beta_{j_*} + \mathcal{O}\left(\frac{1}{\beta_{j_*}} + \frac{\log n}{\sqrt{n}}\right). \quad (112)$$

With the bound (112), we can replace $\langle \mathbf{a}, \mathbf{v} \rangle / (2\beta^2)$ on the righthand side of (11) by $\langle \mathbf{a}, \mathbf{v} \rangle / (2\widehat{\beta}^2)$, which gives an error

$$\left| \frac{\langle \mathbf{a}, \mathbf{v} \rangle}{2\beta_{j_*}^2} - \frac{\langle \mathbf{a}, \mathbf{v} \rangle}{2\widehat{\beta}^2} \right| = \mathcal{O}\left(|\langle \mathbf{a}, \mathbf{v} \rangle| \left(\frac{1}{|\beta_{j_*}|^4} + \frac{\log n}{|\beta_{j_*}|^3 \sqrt{n}} \right)\right).$$

Combining the above discussion together, we can rewrite (11) as

$$\begin{aligned} \langle \mathbf{a}, \widehat{\mathbf{v}} \rangle - \left(1 - \frac{1}{2\widehat{\beta}^2}\right) \langle \mathbf{a}, \mathbf{v} \rangle &= \frac{\langle \mathbf{a}, \boldsymbol{\xi} \rangle - \langle \mathbf{a}, \mathbf{v}_{j_*} \rangle \langle \mathbf{v}_{j_*}, \boldsymbol{\xi} \rangle}{\beta_{j_*}} \\ &+ \mathcal{O}_{\mathbb{P}}\left(\frac{\log n}{\sqrt{n}} \left(\frac{\log n}{\sqrt{n}|\beta_1|}\right)^{k-1} + \frac{\log n}{|\beta_1|^2 \sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta_1|^{3/2} n^{3/4}} + \frac{1}{|\beta_1|^4}\right), \end{aligned} \quad (113)$$

with high probability, where we used that $|\beta_{j_*}| \gtrsim |\beta_1|$.

Again thanks to the definition of $\boldsymbol{\xi}$ in (11) of Theorem 1, i.e. $\boldsymbol{\xi} = \mathbf{Z}[\mathbf{v}^{\otimes(k-1)}]$, is an n -dim vector, with each entry i.i.d. $\mathcal{N}(0, 1/n)$ Gaussian random variable, we see that

$$\langle \mathbf{a}, \boldsymbol{\xi} \rangle - \langle \mathbf{a}, \mathbf{v}_{j_*} \rangle \langle \mathbf{v}_{j_*}, \boldsymbol{\xi} \rangle = \langle \mathbf{a} - \langle \mathbf{a}, \mathbf{v}_{j_*} \rangle \mathbf{v}_{j_*}, \boldsymbol{\xi} \rangle,$$

is a Gaussian random variable, with mean zero and variance

$$\begin{aligned} \mathbb{E}[\langle \mathbf{a} - \langle \mathbf{a}, \mathbf{v}_{j_*} \rangle \mathbf{v}_{j_*}, \boldsymbol{\xi} \rangle^2] &= \frac{1}{n} \|\mathbf{a} - \langle \mathbf{a}, \mathbf{v}_{j_*} \rangle \mathbf{v}_{j_*}\|_2^2 = \frac{1}{n} \langle \mathbf{a}, (\mathbf{I}_n - \mathbf{v}_{j_*} \mathbf{v}_{j_*}^\top) \mathbf{a} \rangle \\ &= \frac{1 + o(1)}{n} \langle \mathbf{a}, (\mathbf{I}_n - \widehat{\mathbf{v}}_{j_*} \widehat{\mathbf{v}}_{j_*}^\top) \mathbf{a} \rangle. \end{aligned}$$

This together with (112), (113) as well as our assumption (13)

$$\frac{\sqrt{n}\widehat{\beta}}{\sqrt{\langle \mathbf{a}, (\mathbf{I}_n - \widehat{\mathbf{v}}\widehat{\mathbf{v}}^\top) \mathbf{a} \rangle}} \left[\left(1 - \frac{1}{2\widehat{\beta}^2}\right)^{-1} \langle \mathbf{a}, \widehat{\mathbf{v}} \rangle - \langle \mathbf{a}, \mathbf{v}_{j_*} \rangle \right] \stackrel{d}{\rightarrow} \mathcal{N}(0, 1). \quad (114)$$

Under the same assumption, we have similar results for Cases 2, 3, 4, by simply changing $(\beta_{j_*}, \mathbf{v}_{j_*})$ in the righthand side of (4) and (5) to the corresponding expression. \blacksquare

Proof [Proof of Corollary 10] For $k \geq 3$ and $|\beta_1| \geq n^{(k-2)/2+\varepsilon}$, the assumption 13 holds trivially. The claim (18) follows from (14). For (19), we recall that in (17), $\boldsymbol{\xi} = \mathbf{Z}[\mathbf{v}_i^{\otimes(k-1)}]$, is an n -dim vector, with each entry i.i.d. $\mathcal{N}(0, 1/n)$ Gaussian random variable. We see that

$$\langle \boldsymbol{\xi}, \mathbf{v}_i \rangle \stackrel{d}{=} \mathcal{N}(0, 1/n).$$

Especially with high probability we will have that $|\langle \boldsymbol{\xi}, \mathbf{v} \rangle| \lesssim \log n / \sqrt{n}$. Then we conclude from (17), with high probability it holds

$$\widehat{\beta} = \beta_i + \mathcal{O}\left(\frac{1}{\beta_i} + \frac{\log n}{\sqrt{n}}\right). \quad (115)$$

With the bound (115), we can replace $(k/2 - 1)/\beta_i$ on the righthand side of (17) by $(k/2 - 1)/\widehat{\beta}$, which gives an error

$$\left| \frac{k/2 - 1}{\beta_i} - \frac{k/2 - 1}{\widehat{\beta}} \right| = \mathcal{O}\left(\frac{1}{|\beta_1|^2} + \frac{\log n}{|\beta_1|\sqrt{n}}\right),$$

where we used that $|\beta_i| \gtrsim |\beta_1|$. Combining the above discussion together, we can rewrite (17) as

$$\beta_i = \widehat{\beta} + \frac{k/2 - 1}{\widehat{\beta}} - \langle \boldsymbol{\xi}, \mathbf{v}_i \rangle + \mathcal{O}_{\mathbb{P}}\left(\frac{\log n}{\sqrt{n}} \left(\frac{\log n}{\sqrt{n}|\beta_1|}\right)^{k-1} + \frac{\log n}{|\beta_1|\sqrt{n}} + \frac{(\log n)^{3/2}}{|\beta_1|^{1/2}n^{3/4}} + \frac{1}{|\beta_1|^2}\right). \quad (116)$$

Since $\langle \boldsymbol{\xi}, \mathbf{v}_i \rangle \stackrel{d}{=} \mathcal{N}(0, 1/n)$, and the error term in (116) is much smaller than $1/\sqrt{n}$. We conclude from (116)

$$\sqrt{n} \left(\beta_i - \widehat{\beta} + \frac{k/2 - 1}{\widehat{\beta}} \right) \xrightarrow{d} \mathcal{N}(0, 1).$$

This finishes the proof of (19). ■

Proof [Proof of Corollary 11] Given the significance level α , the asymptotic confidence intervals in Corollary 11 can be calculated from Corollary 10 by bounding the absolute values of the left hand sides of (18) and (19) at z_α . ■

References

- Emmanuel Abbe, Jianqing Fan, Kaizheng Wang, Yiqiao Zhong, et al. Entrywise eigenvector analysis of random matrices with low expected rank. *Annals of Statistics*, 48(3):1452–1474, 2020.
- Animashree Anandkumar, Rong Ge, Daniel Hsu, and Sham M Kakade. A tensor approach to learning mixed membership community models. *The Journal of Machine Learning Research*, 15(1):2239–2312, 2014a.

- 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, 2014b.
- Gerard Ben Arous, Reza Gheissari, Aukosh Jagannath, et al. Algorithmic thresholds for tensor pca. *Annals of Probability*, 48(4):2052–2087, 2020.
- Arnab Auddy and Ming Yuan. Perturbation bounds for orthogonally decomposable tensors and their applications in high dimensional data analysis. *arXiv preprint arXiv:2007.09024*, 2020.
- Zhidong Bai and Jianfeng Yao. On sample eigenvalues in a generalized spiked population model. *Journal of Multivariate Analysis*, 106:167–177, 2012.
- Jinho Baik and Jack W Silverstein. Eigenvalues of large sample covariance matrices of spiked population models. *Journal of multivariate analysis*, 97(6):1382–1408, 2006.
- Jinho Baik, Gérard Ben Arous, Sandrine Péché, et al. Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices. *The Annals of Probability*, 33(5):1643–1697, 2005.
- Mikhail Belkin, Luis Rademacher, and James Voss. Eigenvectors of orthogonally decomposable functions. *SIAM Journal on Computing*, 47(2):547–615, 2018.
- Florent Benaych-Georges and Raj Rao Nadakuditi. The singular values and vectors of low rank perturbations of large rectangular random matrices. *Journal of Multivariate Analysis*, 111:120–135, 2012.
- Aharon Birnbaum, Iain M Johnstone, Boaz Nadler, and Debashis Paul. Minimax bounds for sparse pca with noisy high-dimensional data. *Annals of statistics*, 41(3):1055, 2013.
- T Tony Cai, Zongming Ma, Yihong Wu, et al. Sparse pca: Optimal rates and adaptive estimation. *The Annals of Statistics*, 41(6):3074–3110, 2013.
- Tony Cai, Zongming Ma, and Yihong Wu. Optimal estimation and rank detection for sparse spiked covariance matrices. *Probability theory and related fields*, 161(3-4):781–815, 2015.
- Jie Chen and Yousef Saad. On the tensor svd and the optimal low rank orthogonal approximation of tensors. *SIAM journal on Matrix Analysis and Applications*, 30(4):1709–1734, 2009.
- Wei-Kuo Chen, Madeline Handschy, and Gilad Lerman. Phase transition in random tensors with multiple spikes. *arXiv preprint arXiv:1809.06790*, 2018a.
- Wei-Kuo Chen et al. Phase transition in the spiked random tensor with rademacher prior. *The Annals of Statistics*, 47(5):2734–2756, 2019.
- Yuxin Chen, Chen Cheng, and Jianqing Fan. Asymmetry helps: Eigenvalue and eigenvector analyses of asymmetrically perturbed low-rank matrices. *arXiv preprint arXiv:1811.12804*, 2018b.

- Chen Cheng, Yuting Wei, and Yuxin Chen. Inference for linear forms of eigenvectors under minimal eigenvalue separation: Asymmetry and heteroscedasticity. *arXiv preprint arXiv:2001.04620*, 2020.
- Andrzej Cichocki, Danilo Mandic, Lieven De Lathauwer, Guoxu Zhou, Qibin Zhao, Cesar Caiafa, and Huy Anh Phan. Tensor decompositions for signal processing applications: From two-way to multiway component analysis. *IEEE signal processing magazine*, 32(2): 145–163, 2015.
- Pierre Comon. Tensors: a brief introduction. *IEEE Signal Processing Magazine*, 31(3): 44–53, 2014.
- David L Donoho, Matan Gavish, and Iain M Johnstone. Optimal shrinkage of eigenvalues in the spiked covariance model. *Annals of statistics*, 46(4):1742, 2018.
- Olivier Duchenne, Francis Bach, In-So Kweon, and Jean Ponce. A tensor-based algorithm for high-order graph matching. *IEEE transactions on pattern analysis and machine intelligence*, 33(12):2383–2395, 2011.
- Noureddine El Karoui et al. Spectrum estimation for large dimensional covariance matrices using random matrix theory. *The Annals of Statistics*, 36(6):2757–2790, 2008.
- Evgeny Frolov and Ivan Oseledets. Tensor methods and recommender systems. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 7(3):e1201, 2017.
- Wolfgang Hackbusch. *Tensor spaces and numerical tensor calculus*, volume 42. Springer, 2012.
- Rungang Han, Rebecca Willett, and Anru Zhang. An optimal statistical and computational framework for generalized tensor estimation. *arXiv preprint arXiv:2002.11255*, 2020.
- Christopher J Hillar and Lek-Heng Lim. Most tensor problems are np-hard. *Journal of the ACM (JACM)*, 60(6):1–39, 2013.
- Samuel B Hopkins, Jonathan Shi, and David Steurer. Tensor principal component analysis via sum-of-square proofs. In *Conference on Learning Theory*, pages 956–1006, 2015.
- Samuel B Hopkins, Tselil Schramm, Jonathan Shi, and David Steurer. Fast spectral algorithms from sum-of-squares proofs: tensor decomposition and planted sparse vectors. In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pages 178–191, 2016.
- 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, 2013.
- Aukosh Jagannath, Patrick Lopatto, and Leo Miolane. Statistical thresholds for tensor pca. *arXiv preprint arXiv:1812.03403*, 2018.

- Bing-Yi Jing, Ting Li, Zhongyuan Lyu, and Dong Xia. Community detection on mixture multi-layer networks via regularized tensor decomposition. *arXiv preprint arXiv:2002.04457*, 2020.
- Iain M Johnstone. On the distribution of the largest eigenvalue in principal components analysis. *Annals of statistics*, pages 295–327, 2001.
- Iain M Johnstone and Arthur Yu Lu. On consistency and sparsity for principal components analysis in high dimensions. *Journal of the American Statistical Association*, 104(486): 682–693, 2009.
- Iain M Johnstone and Debashis Paul. PCA in high dimensions: An orientation. *Proceedings of the IEEE*, 106(8):1277–1292, 2018.
- Alexandros Karatzoglou, Xavier Amatriain, Linas Baltrunas, and Nuria Oliver. Multiverse recommendation: n-dimensional tensor factorization for context-aware collaborative filtering. In *Proceedings of the fourth ACM conference on Recommender systems*, pages 79–86, 2010.
- Chiheon Kim, Afonso S Bandeira, and Michel X Goemans. Community detection in hypergraphs, spiked tensor models, and sum-of-squares. In *2017 International Conference on Sampling Theory and Applications (SampTA)*, pages 124–128. IEEE, 2017.
- Tamara G Kolda. Orthogonal tensor decompositions. *SIAM Journal on Matrix Analysis and Applications*, 23(1):243–255, 2001.
- Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. *SIAM review*, 51(3):455–500, 2009.
- Olivier Ledoit, Michael Wolf, et al. Nonlinear shrinkage estimation of large-dimensional covariance matrices. *The Annals of Statistics*, 40(2):1024–1060, 2012.
- Thibault Lesieur, Léo Miolane, Marc Lelarge, Florent Krzakala, and Lenka Zdeborová. Statistical and computational phase transitions in spiked tensor estimation. In *2017 IEEE International Symposium on Information Theory (ISIT)*, pages 511–515. IEEE, 2017.
- Yuetian Luo and Anru R Zhang. Open problem: Average-case hardness of hypergraphic planted clique detection. In *Conference on Learning Theory*, pages 3852–3856. PMLR, 2020a.
- Yuetian Luo and Anru R Zhang. Tensor clustering with planted structures: Statistical optimality and computational limits. *arXiv preprint arXiv:2005.10743*, 2020b.
- Yuetian Luo, Garvesh Raskutti, Ming Yuan, and Anru R Zhang. A sharp blockwise tensor perturbation bound for orthogonal iteration. *arXiv preprint arXiv:2008.02437*, 2020.
- Zongming Ma et al. Sparse principal component analysis and iterative thresholding. *The Annals of Statistics*, 41(2):772–801, 2013.

- Sean O’Rourke, Van Vu, and Ke Wang. Random perturbation of low rank matrices: Improving classical bounds. *Linear Algebra and its Applications*, 540:26–59, 2018.
- Debashis Paul. Asymptotics of sample eigenstructure for a large dimensional spiked covariance model. *Statistica Sinica*, pages 1617–1642, 2007.
- Sandrine Péché. The largest eigenvalue of small rank perturbations of hermitian random matrices. *Probability Theory and Related Fields*, 134(1):127–173, 2006.
- Amelia Perry, Alexander S Wein, Afonso S Bandeira, et al. Statistical limits of spiked tensor models. In *Annales de l’Institut Henri Poincaré, Probabilités et Statistiques*, volume 56, pages 230–264. Institut Henri Poincaré, 2020.
- Steffen Rendle and Lars Schmidt-Thieme. Pairwise interaction tensor factorization for personalized tag recommendation. In *Proceedings of the third ACM international conference on Web search and data mining*, pages 81–90, 2010.
- Emile Richard and Andrea Montanari. A statistical model for tensor pca. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2897–2905. Curran Associates, Inc., 2014. URL <http://papers.nips.cc/paper/5616-a-statistical-model-for-tensor-pca.pdf>.
- Elina Robeva. Orthogonal decomposition of symmetric tensors. *SIAM Journal on Matrix Analysis and Applications*, 37(1):86–102, 2016.
- Nicholas D Sidiropoulos, Lieven De Lathauwer, Xiao Fu, Kejun Huang, Evangelos E Papalexakis, and Christos Faloutsos. Tensor decomposition for signal processing and machine learning. *IEEE Transactions on Signal Processing*, 65(13):3551–3582, 2017.
- Erez Simony, Christopher J Honey, Janice Chen, Olga Lositsky, Yaara Yeshurun, Ami Wiesel, and Uri Hasson. Dynamic reconfiguration of the default mode network during narrative comprehension. *Nature communications*, 7:12141, 2016.
- Van Vu. Singular vectors under random perturbation. *Random Structures & Algorithms*, 39(4):526–538, 2011.
- Vincent Q Vu, Jing Lei, et al. Minimax sparse principal subspace estimation in high dimensions. *The Annals of Statistics*, 41(6):2905–2947, 2013.
- Dong Xia, Anru R Zhang, and Yuchen Zhou. Inference for low-rank tensors—no need to debias. *arXiv preprint arXiv:2012.14844*, 2020.
- Anru Zhang and Dong Xia. Tensor svd: Statistical and computational limits. *IEEE Transactions on Information Theory*, 64(11):7311–7338, 2018.
- Anru Zhang, T Tony Cai, and Yihong Wu. Heteroskedastic pca: Algorithm, optimality, and applications. *arXiv preprint arXiv:1810.08316*, 2018.
- Yiqiao Zhong. Eigenvector under random perturbation: A nonasymptotic rayleigh-schrödinger theory. *arXiv preprint arXiv:1702.00139*, 2017.

Hua Zhou, Lexin Li, and Hongtu Zhu. Tensor regression with applications in neuroimaging data analysis. *Journal of the American Statistical Association*, 108(502):540–552, 2013.