

# Stochastic Approximation for Online Tensorial Independent Component Analysis

Chris Junchi Li  
 Michael I. Jordan

JUNCHILI@BERKELEY.EDU  
 JORDAN@CS.BERKELEY.EDU

*Department of EECS & Department of Statistics, University of California, Berkeley*

**Editors:** Mikhail Belkin and Samory Kpotufe

## Abstract

Independent component analysis (ICA) has been a popular dimension reduction tool in statistical machine learning and signal processing. In this paper, we present a convergence analysis for an online tensorial ICA algorithm, by viewing the problem as a nonconvex stochastic approximation problem. For estimating one component, we provide a dynamics-based analysis to prove that our online tensorial ICA algorithm with a specific choice of stepsize achieves a sharp finite-sample error bound. In particular, under a mild assumption on the data-generating distribution and a scaling condition such that  $d^4/T$  is sufficiently small up to a polylogarithmic factor of data dimension  $d$  and sample size  $T$ , a sharp finite-sample error bound of  $\tilde{O}(\sqrt{d/T})$  can be obtained.

**Keywords:** Independent component analysis, tensor decomposition, non-Gaussianity, finite-sample error bound, online learning

## 1. Introduction

Independent Component Analysis (ICA) is a widely used dimension reduction method with diverse applications in the fields of statistical machine learning and signal processing (Hyvärinen et al., 2001; Stone, 2004; Samworth and Yuan, 2012). Let the data vector be modeled as  $\mathbf{X} = \mathbf{A}\mathbf{Z}$ , where  $\mathbf{A} \equiv (\mathbf{a}_1, \dots, \mathbf{a}_d) \in \mathbb{R}^{d \times d}$  is a full-rank mixing matrix whose columns are orthogonal components in  $\mathbb{R}^d$ , and  $\mathbf{Z} \in \mathbb{R}^d$  is a non-Gaussian latent random vector consisting of independent entries  $\mathbf{Z} = (Z_1, \dots, Z_d)^\top$ . The goal of ICA is to recover one or multiple columns among  $(\mathbf{a}_1, \dots, \mathbf{a}_d)$  from independent observations of  $\mathbf{X} = (X_1, \dots, X_d)^\top$ . Following standard practice, we assume that the random vector  $\mathbf{X}$  has been whitened in the sense that it has zero mean and an identity covariance matrix (Hyvärinen et al., 2001), and we focus on the case where the distributions of  $Z_1, \dots, Z_d$  share a fourth moment  $\mu_4 \neq 3$  (i.e., they are of identical kurtosis,  $\mu_4 - 3$ ). These assumptions restrict the search of mixing matrix  $\mathbf{A}$  to the space of orthogonal matrices and guarantee its identifiability up to signed permutations of its columns  $(\mathbf{a}_1, \dots, \mathbf{a}_d)$  (Comon, 1994; Frieze et al., 1996; Hyvärinen et al., 2001).

In this paper, we study a stochastic algorithm that estimates independent components for streaming data. Such an algorithm processes and discards one or a small batch of data observations at each iterate and enjoys reduced storage complexity. To begin with, we cast the tensorial ICA as the problem of optimizing an objective function based on the fourth-order cumulant tensor over the unit sphere  $\mathcal{D}_1 \equiv \{\mathbf{u} \in \mathbb{R}^d : \|\mathbf{u}\| = 1\}$ , where  $\|\cdot\|$  denotes the Euclidean norm. The optimization problem is as follows:

$$\begin{aligned} \min \quad & -\text{sign}(\mu_4 - 3)\mathbb{E}(\mathbf{u}^\top \mathbf{X})^4 \\ \text{subject to} \quad & \mathbf{u} \in \mathcal{D}_1. \end{aligned} \tag{1.1}$$

Such an objective function is referred to as a non-Gaussianity contrast function (Hyvärinen et al., 2001). The formulation (1.1) of ICA via its fourth-order cumulant tensor was initiated by Comon (1994) and Frieze et al. (1996) in the batch case. The landscape of objective (1.1) is featured by its nonconvexity, in the sense that it presents local minimizers  $a_1, \dots, a_d$  and exponentially many unstable stationary points in terms of (Ge et al., 2015; Li et al., 2016; Sun et al., 2015). Here, by unstable stationary points we refer to those points with zero gradient and at least one negative Hessian direction, and hence includes the collection of saddle points and local maximizers. Algorithms that find local minimizers of (1.1) allow us to estimate the columns of the mixing matrix  $A$ . We analyze and discuss the following stochastic approximation method, which we refer to as online tensorial ICA (Ge et al., 2015; Li et al., 2016; Wang and Lu, 2017). Initialize a unit vector  $u^{(0)}$  appropriately, at step  $t = 1; 2; \dots; T$  the algorithm processes an observation  $x^{(t)}$  by performing the following update:

$$u^{(t)} = \frac{1}{\| \cdot \|} \left( u^{(t-1)} + \eta^{(t)} \text{sign}(\langle \cdot, \cdot \rangle) (u^{(t-1)})^\top X^{(t)3} X^{(t)} \right); \quad (1.2)$$

where  $\eta^{(t)}$  is a positive stepsize, and the operator  $[\cdot] := \frac{\cdot}{\|\cdot\|}$  projects a nonzero vector onto the unit sphere  $\mathbb{D}_1$  centered at the origin.

Due to the nonconvexity of the optimization problem (1.1), a key issue that arises in analyzing the iteration (1.2) is the avoidance of unstable stationary points including saddle points and local maximizers. In earlier work, Ge et al. (2015) introduced an additional artificial noise injection step to the algorithm and developed a hit-and-escape convergence analysis, obtaining a polynomial convergence rate for the online tensorial ICA problem and beyond (for a more general class of nonconvex optimization problems). In contrast, we present a dynamics-based analysis that requires no noise injection step. We show that a uniform initialization on the unit sphere along with mild scaling conditions, is sufficient to ensure that the algorithm iterate enters a basin of attraction and finds a local (and also global) minimizer with high probability.

**Overview of the Main Result and Contributions** Our main result states that for estimating a single independent component, the iteration (1.2) with uniform initialization and carefully chosen stepsize  $\eta^{(t)}$  enters the basin of attraction of a uniformly drawn independent component and achieves a convergence rate of  $\frac{1}{\sqrt{d}} \frac{1}{T}$  in terms of angle with respect to an independent component, up to a polylogarithmic factor of  $d$  and  $T$ . Such a convergence result holds under mild distributional assumptions and scaling conditions that involve the data dimension and sample size  $T$ . Informally, this result is stated as follows.

**Theorem 1 (Informal version of Corollary 8 in x3)** Let the initial  $u^{(0)}$  be sampled uniformly from the unit sphere  $\mathbb{D}_1$ , and let appropriate light-tailed distributional assumptions hold. Then, for any fixed positive  $\epsilon \in (0; 1/5]$  satisfying the scaling condition:

$$d \leq \frac{2^p}{2e} \log^{-1} \frac{1}{\epsilon} + 1; \quad C_{1;T} \log^8(C_{1;T}^{-1} d T) \leq \frac{B^8}{j^4 3j^2} \frac{d^4 \log^2 T}{T} \leq \frac{\epsilon^2}{\log^2 \frac{1}{\epsilon}};$$

1. In contrast, the rank-one eigenvalue problem presents local minimizers and  $2d - 2$  unstable stationary points in the typical setting of distinct eigenvalues.

---

Algorithm 1 Online Tensorial ICA

---

Initialize  $u^{(0)}$  and select stepsize  $\epsilon^{(t)}$  appropriately (to elaborate later)  
 for  $t = 1; 2; \dots$  do  
     Draw one observation  $X^{(t)}$  from streaming data, and update iteration  $t$  via

$$u^{(t)} = \Pi_{D_1} u^{(t-1)} + \epsilon^{(t)} \text{sign}(\langle u^{(t-1)}, X^{(t)} \rangle) X^{(t)} \quad (1.3)$$

where  $\Pi_{D_1} g = \frac{g}{\|g\|}$  denotes the projection operator onto the unit sphere centered at the origin  $D_1$   
 end for

---

there exists a uniformly distributed random variable  $\epsilon \in [d]$  such that with probability at least  $1 - \delta$ , iteration  $u^{(t)}$  of (1.2) with appropriate choice of stepsize  $\epsilon^{(t)}$  satisfies

$$\tan \angle(u^{(T)}; a_\epsilon) \leq C_{1;T} \log^{5/2}(C_{1;T} \epsilon^{-1} d) \frac{B^4}{j} \frac{d \log^2 T}{T};$$

where  $C_{1;T}$  is a positive, absolute constant.

To the best of our knowledge, this provides the first rigorous analysis of an online tensorial ICA algorithm that achieves a  $\tilde{O}(\frac{1}{\sqrt{d}})$  nite-sample convergence rate, under mild distributional assumptions and scaling conditions. The contributions of this work lie in several aspects:

- (i) Partly adapting from the analysis of online principal component estimation in Li et al. (2018), we provide a per-coordinate analysis of the warmly-initialized online tensorial ICA algorithm that achieves a sharp  $\tilde{O}(\frac{1}{\sqrt{d}})$  convergence rate [Theorem 4].
- (ii) In contradistinction to existing saddle-point-escaping analysis in Ge et al. (2015); Jin et al. (2017); Li et al. (2018), we developed a novel coordinate-pair analysis of the uniformly-initialized online tensorial ICA algorithm based on our dynamics-based characterization [Theorem 7].
- (iii) We combine these two analyses to conclude that a two-stage training procedure provides a nite-sample error bound  $\tilde{O}(\frac{1}{\sqrt{d}})$  for the uniform initialization case under a mild assumption on the data-generating distribution and also a scaling condition  $\epsilon = \tilde{O}(\frac{1}{\sqrt{d}})$  being sufficiently small, up to a polylogarithmic factor of  $d; T$  [Theorem 1, presented formally in Corollary 8].

Organization The rest of this paper is organized as follows. Section 2 presents our main convergence results and nite-sample error bounds for the warm initialization case for estimating one single component. Section 3 presents the corresponding results for the uniform initialization case. Section 4 discusses additional related literature. Section 5 summarizes our results. Limited by space we relegate to the Appendix the proofs of main results, all secondary lemmas and technical results along with preliminary simulation results.

## 2. Estimating a Single Component: Warm Initialization Case

For the purpose of estimating a single independent component, we introduce our settings and assumptions for tensorial ICA and its stochastic approximation algorithm (1.2), formally stated in

**Algorithm 1.** Let the dimension  $d \geq 2$ , let  $X$  be the data vector of which  $X^{(1)}; X^{(2)}; \dots; X^{(n)} \in \mathbb{R}^d$  are independent draws, and assume the following for the distribution of  $X$ :

**Assumption 2 (Data vector distribution)** Let  $X = AZ$ , where  $A \in \mathbb{R}^{d \times d}$  is an orthogonal matrix with  $\|A\|_F = 1$  and  $Z \in \mathbb{R}^d$  is a random vector satisfying

- (i) The  $Z_i; i = 1, \dots, d$  are independent with identical  $j$ -th-moment for  $j = 1, 2, 4$ , denoted as  $\mathbb{E} Z_i^j$ ;
- (ii) The  $\mu_1 = \mathbb{E} Z_i = 0, \mu_2 = \mathbb{E} Z_i^2 = 1, \mu_4 = \mathbb{E} Z_i^4 \in (3, 9)$ ;
- (iii) For all  $i \in [d]$ ,  $Z_i$  has an Orlicz- $\psi_2$  norm bounded by  $\psi_2(Z_i) \leq B$ .

Note that Assumption 2(i) requires the distribution for all independent components to admit identical first, second and fourth moments. As indicated in Assumption 2(ii), the data vectors are assumed to be whitened first in the sense that  $\mu_2 = 1$ . The sign of our excess kurtosis  $\mu_4 - 3$  determines the direction of stochastic gradient update, and, as will be seen later, the magnitude of the excess kurtosis  $\mu_4 - 3$  plays an important role in our convergence analysis. Assumption 2(iii) generalizes the boundedness assumption  $\psi_2(Z_i) \leq O(B)$  made in recent work in order to cover Gaussian mixtures and Bernoulli-Gaussians, which are typical application cases for the tensorial ICA estimation problem. Note that we include a factor of  $3/8$  in the Orlicz- $\psi_2$  norm (sub-Gaussian parameter) purely for notational simplicity in our analysis.

We target to study the convergence of tensorial ICA under certain initialization conditions. For each initialization condition, we first analyze the convergence result for any fixed, plausible step-sizes, and then (by choosing the stepsize according to the number of observations) obtain the finite-sample error bound. We focus in this section the warm initialization condition as any satisfying, for some integer  $2 \leq [d]$ ,

$$u^{(0)} \in D_1 \text{ and } \tan \angle(u^{(0)}; a_i) \leq \frac{1}{3} \tag{2.1}$$

For any fixed  $\epsilon > 0$ , we define a rescaled time variable  $T_\epsilon$  as

$$T_\epsilon := \frac{2}{\epsilon} \frac{\log \frac{1 - \frac{4}{3} \epsilon^3}{B^8}}{\log \frac{1 - \frac{4}{3} \epsilon^3}{4} - \frac{1}{3} \epsilon^3} \tag{2.2}$$

Then under warm initialization condition (2.1), we have the following convergence lemma.

**Lemma 3 (Convergence Result with Warm Initialization)** Let the dimension  $d \geq 2$ , let Assumption 2 hold, and let initialization  $u^{(0)}$  satisfy condition (2.1) for some integer  $2 \leq [d]$ . Then, for any

2. A random variable  $Z$  with mean 0 is light-tailed if there is a positive number  $\epsilon$  such that  $\mathbb{E} \exp(\epsilon Z) \leq \exp(\epsilon^2/2)$  for all  $\epsilon \in \mathbb{R}$ . The smallest possible  $\epsilon > 0$  is referred to as the Orlicz- $\psi_2$  norm or sub-Gaussian parameter (Wainwright, 2019). Readers shall be warned that the term “sub-Gaussian” here indicates the light-tailed condition and should be distinguished from  $\mu_4 < 3$  (resp.  $\mu_4 > 3$ ) of sub-Gaussianity (resp. super-Gaussianity) in ICA (Stone, 2004).
3. When the excess kurtosis  $\mu_4 - 3$  is equal to zero, for instance when  $Z$  follows an i.i.d. standard normal distribution, the matrix  $A$  is non-identifiable in our framework. Non-gaussian independent components with  $\mu_4 > 3$  can be studied via higher-order tensor decomposition with a different contrast function, but is beyond the scope of this paper.

fixed positive numbers  $\bar{\epsilon}$  and  $\beta \in (0, 1]$  satisfying the scaling condition,

$$C_{3,L} \log^8(T; \beta^{-1}) \frac{B^8}{j^4 3^j} d \log \frac{j^4 3^j}{B^8} \beta^{-1} \frac{1}{\beta+1}; \quad < \min \left\{ \frac{1}{j^4 3^j}; \frac{j^4 3^j}{B^8} e^{-1} \right\}; \quad (2.3)$$

there exists an event  $\mathcal{H}_{3,L}$  with  $P(\mathcal{H}_{3,L}) \geq 1 - 6\beta + 12\beta + \frac{5184}{\log^5 \beta^{-1}} \beta^{-1} d$ ; such that on  $\mathcal{H}_{3,L}$ , iteration  $u^{(t)}$  of Algorithm 1 satisfies

$$\begin{aligned} \tan \backslash u^{(t)}; a_i & \leq \tan \backslash u^{(0)}; a_i \beta^{-1} \frac{j^4 3^j}{3^j} \beta^{-t} \\ & + \frac{\beta}{\beta+1} C_{3,L} \log^{5=2} \beta^{-1} \frac{B^4}{j^4 3^{j=2}} \beta^{-1} d \log \frac{j^4 3^j}{B^8} \beta^{-1}; \end{aligned} \quad (2.4)$$

for all  $t \in [0; T; ]$ , where  $C_{3,L}$  and  $C_{3,L}$  are positive, absolute constants.

Lemma 3 provides, under the warm initialization assumption (2.1) and scaling condition  $\mathcal{O}(d^{-1})$ , an upper bound for  $\tan \backslash u^{(t)}; a_i$  which is the sum of two terms: the first term on the right hand side of (2.4) decays geometrically from  $\tan \backslash u^{(0)}; a_i$  at rate  $\frac{j^4 3^j}{3^j} = 3$ , and the second term  $\frac{\beta}{\beta+1} C_{3,L} \log^{5=2} \beta^{-1} \frac{B^4}{j^4 3^{j=2}} \beta^{-1} d \log \frac{j^4 3^j}{B^8} \beta^{-1}$  is induced by the noise. To balance these two terms, when we know in advance the sample size of online data satisfying some scaling condition  $\beta = \beta(d)$ , we choose a constant stepsize  $\beta = \beta(\log T = T)$  and obtain a finite-sample error bound:

**Theorem 4 (Finite-Sample Error Bound with Warm Initialization)** Let the dimension  $d \geq 2$ , let Assumption 2 hold, and let initialization  $u^{(0)}$  satisfy condition (2.1) for some integer  $\beta \in [d]$ . For sample size  $T$  set the stepsize as

$$\beta(T) = \frac{9 \log \frac{2j^4 3^{j^2} T}{9B^8}}{2j^4 3^j T}; \quad (2.5)$$

Then, for any fixed positive numbers  $\bar{\epsilon} \in (0, 1]$  and  $\beta \in (0, 1]$  satisfying the scaling condition

$$C_{4,T} \log^8(C_{4,T}^0 \beta^{-1} d T) \frac{d \log^2 T}{T} \frac{j^4 3^{j^2}}{B^8}; \quad (2.6)$$

there exists an event  $\mathcal{H}_{4,T}$  with  $P(\mathcal{H}_{4,T}) \geq 1 - \bar{\epsilon}$  such that on  $\mathcal{H}_{4,T}$ , iteration  $u^{(t)}$  of Algorithm 1 satisfies

$$\tan \backslash u^{(T)}; a_i \leq C_{4,T} \log^{5=2}(C_{4,T}^0 \beta^{-1} d) \frac{B^4}{j^4 3^j} \frac{d \log^2 T}{T};$$

where  $C_{4,T}; C_{4,T}; C_{4,T}^0$  are positive, absolute constants.

Theorem 4 indicates a  $\frac{\beta}{\beta+1} \frac{1}{d=T}$  finite-sample error bound for online tensorial ICA when it is warmly initialized in the sense that (2.1) holds for  $u^{(0)}$  for some integer  $\beta \in [d]$ . In the rest of this section, we prove Lemma 3 and Theorem 4 for the warm initialization case. The section is organized as follows: 2.1 analyzes our algorithm and provide a key lemma [Lemma 5] when it is warmly initialized. Proof of the key Lemma 5 and all proofs of secondary lemmas are deferred to xA.1 and xB in Appendix.

2.1. Key Lemma in the Warm Initialization Analysis

To simplify our problem, we set  $2 \leq d$  as the integer such that  $\mathbf{u}^{(0)}$  satisfies condition (2.1) and define the rotated iteration  $\mathbf{v}^{(t)} \in \mathbb{R}^d$  as

$$\mathbf{v}^{(t)} = \mathbf{P} \mathbf{A}^T \mathbf{u}^{(t)}; \tag{2.7}$$

where  $\mathbf{P} \in \mathbb{R}^{d \times d}$  is the permutation matrix corresponding to the cycle; i.e.,  $P(i; 1) = P(1; i) = 1$ ,  $P(j; j) = 1$  for  $j \in \{2, \dots, d\}$  and all other elements being zero. Such a matrix, as an operator, maps the component vectors to coordinate vectors and ensures that the closest independent components pair at initialization and (with high probability) at convergence. Furthermore, transforming to rotated observations  $\mathbf{Y}^{(t)} = \mathbf{P} \mathbf{A}^T \mathbf{X}^{(t)}$  allows us to equivalently translate our online tensorial ICA iteration (1.3) into an analogous form:

$$\mathbf{v}^{(t)} = \frac{1}{\|\mathbf{v}^{(t-1)}\|} \mathbf{v}^{(t-1)} + \frac{\text{sign}(\langle \mathbf{v}^{(t-1)}, \mathbf{Y}^{(t)} \rangle)}{\|\mathbf{v}^{(t-1)}\|^3} \mathbf{Y}^{(t)}; \tag{2.8}$$

It is easy to verify that the rotated iterations  $\mathbf{v}^{(t)} \in \mathbb{R}^d$  and  $\mathbf{u}^{(t)} \in \mathbb{R}^d$  satisfy  $\mathbf{a}_i^T \mathbf{u}^{(t)} = \mathbf{e}_1^T \mathbf{v}^{(t)}$ , and hence for all  $t \geq 0$

$$\tan \angle(\mathbf{v}^{(t)}; \mathbf{e}_1) = \frac{\frac{1}{\|\mathbf{v}^{(t)}\|} \langle \mathbf{v}^{(t)}, \mathbf{e}_1 \rangle}{\|\mathbf{v}^{(t)}\|} = \frac{1}{\|\mathbf{v}^{(t)}\|} \frac{\langle \mathbf{v}^{(t)}, \mathbf{e}_1 \rangle}{\|\mathbf{v}^{(t)}\|} = \frac{1}{\|\mathbf{v}^{(t)}\|} \frac{\langle \mathbf{v}^{(t)}, \mathbf{e}_1 \rangle}{\|\mathbf{v}^{(t)}\|} = \tan \angle(\mathbf{u}^{(t)}; \mathbf{a}_i); \tag{2.9}$$

Now, we let the warm initialization region be

$$\mathcal{D}_{\text{warm}} = \{ \mathbf{v} \in \mathbb{R}^d : \mathbf{v}_1^2 \geq \frac{3}{4} \} = \{ \mathbf{v} \in \mathbb{R}^d : \tan \angle(\mathbf{v}; \mathbf{e}_1) \leq \frac{1}{\sqrt{3}} \}; \tag{2.10}$$

Note that the warm initialization condition in (2.1) is equivalent to  $\mathbf{v}^{(0)} \in \mathcal{D}_{\text{warm}}$ . Analogous to the warm initialization studied in the setting of principal component estimation (Li et al., 2018), the iteration we study in our warmly initialized online tensorial ICA is

$$\mathbf{U}_k^{(t)} = \frac{\mathbf{v}_k^{(t)}}{\|\mathbf{v}_k^{(t)}\|}; \tag{2.11}$$

To prevent iteration  $\mathbf{U}_k^{(t)}$  from diverging from the warm initialization region, we also define a larger warm-auxiliary region as  $\mathcal{D}_{\text{warm-aux}} = \{ \mathbf{v} \in \mathbb{R}^d : \mathbf{v}_1^2 \geq \frac{2}{3} \} = \{ \mathbf{v} \in \mathbb{R}^d : \tan \angle(\mathbf{v}; \mathbf{e}_1) \leq \frac{1}{\sqrt{2}} \}$ .

Suppose the process  $\mathbf{v}^{(t)}$  is initialized at  $\mathbf{v}^{(0)} \in \mathcal{D}_{\text{warm}}$ , and we define the first time  $\mathbf{v}^{(t)}$  exits the warm-auxiliary region as

$$T_x = \inf_{t \geq 1} \{ \mathbf{v}^{(t)} \notin \mathcal{D}_{\text{warm-aux}}^c \}; \tag{2.12}$$

We state the key lemma as

**Lemma 5** Let the settings in Lemma 3 hold, and let  $x$  the coordinate  $2 \leq x \leq d$  and value  $\epsilon > 0$ . Then for any fixed positives  $\epsilon, \delta$  satisfying the scaling condition (2.3) along with the warm

initialization condition  $v^{(0)} \in D_{\text{warm}}$ , there exists an event  $H_{k;5;L}$  satisfying  $P(H_{k;5;L}) \geq 1 - 6^{-6} - 12 + \frac{5184}{\log^5 - 1}$ ; such that on event  $H_{k;5;L}$  the following holds:

$$\sup_{t \in [T_{;0;5}^0, T_x^0]} \|U_k^{(t)} - U_k^{(0)}\|_F \leq \frac{1}{j-4} \frac{3j}{2d} \sum_{s=0}^t (v_1^{(s)})^2 + (v_k^{(s)})^2 \leq 2C_{5;L} \log^{5-2} \frac{1}{B^4} (T_{;0;5}^0)^{1-2}; \quad (2.13)$$

where  $C_{5;L}$  is a positive, absolute constant.

This lemma shows that for each coordinate  $u_k^{(t)}$  [2; d], with high probability, the dynamics  $u_k^{(t)}$  is tightly controlled within a deterministic vessel whose center converges to zero at least exponentially fast. As we will later see in the proofs of Lemma 3 and Theorem 4, this guarantees that with high probability  $v^{(t)}$  will not exit the warm-auxilliary region  $D_{\text{warm-aux}}$  and will stay within a small neighborhood of  $e_1$  after  $T_{;0;5}^0$  iterates, where the rescaled time  $\tau$  was earlier defined in (2.2).

### 3. Estimating Single Component: Uniform Initialization Case

It is often the case that we are unable to obtain a warm initialization for online tensorial ICA as required in x2, in which case we resort to initializing  $u^{(0)}$  uniformly at random from the unit sphere  $D_1$ . To analyze such a case, we define a new rescaled time variable, parameterized by  $\gamma$  as follows:

$$T_{;0;5}^0 = \frac{2}{6} \frac{\log \frac{j-4}{B^8} \frac{3j}{2d} \frac{1}{3j}}{\log 1 - \frac{1}{2d} \frac{3j}{4}} \frac{7}{7}; \quad (3.1)$$

Then under uniform initialization, we have the following convergence result:

Lemma 6 (Convergence Result with Uniform Initialization) Let the dimension  $d \geq 2$ , let Assumption 2 hold, and let  $u^{(0)}$  be uniformly sampled from the unit sphere  $D_1$ . Then for any fixed positive numbers  $\delta \geq 2$ ;  $\gamma > 0.5$ ;  $\epsilon > 0$ ;  $\epsilon_2 \in (0, e^{-1}]$ ;  $\epsilon_3 \in (0, 1-3]$  satisfying the scaling condition

$$\delta \geq 2^{\frac{p}{2}} \frac{1}{e \log^{-1} + 1}; \quad \epsilon_2 < \min \left\{ \frac{1}{j-4} \frac{3j}{2d}, \frac{j-4}{B^8} e^{-1} \right\}; \quad \text{and} \quad (3.2)$$

$$C_{6;L} \log^8(T_{;0;5}^0) \leq \frac{B^8}{j-4} \frac{3j}{2d} \delta^2 \log^2 d \log \frac{j-4}{B^8} \frac{3j}{2d} \frac{1}{3j} \frac{2}{(\delta+1) \log^2 - 1};$$

there exists a uniformly distributed random variable  $u \in \mathbb{R}^d$  and an event  $H_{6;L}$  with

$$P(H_{6;L}) \geq 1 - 6^{-6} - 27 + \frac{10368}{\log^5 - 1} \delta^{-3};$$

such that on  $H_{6;L}$ , iteration  $u^{(t)}$  of Algorithm 1 satisfies for  $t \in [T_{;0;5}^0, T_x^0]$

$$\|u^{(t)} - a\|_2 \leq \frac{p}{d} \frac{1}{2d} \frac{3j}{4} \sum_{s=0}^t \tau^{T_{;0;5}^0} + \frac{p}{\delta+1} C_{6;L} \log^{5-2} \frac{1}{B^4} \frac{1}{d^3 \log \frac{j-4}{B^8} \frac{3j}{2d} \frac{1}{3j}}; \quad (3.3)$$

where  $C_{6;L}$ ;  $C_{6;L}$  are positive, absolute constants.

Under the uniform initialization assumption, Lemma 6 shows that the term  $\tan \setminus u^{(t)}; a_i$  can be upper bounded by the sum of two summands, with the first term geometrically decaying from  $\frac{1}{d}$  at rate  $\frac{1}{j-4-3j} = (2d)^{-j}$ , and the second being the noise-induced error  $\frac{1}{d^3}$ . The key idea behind the proof is that the scaling condition (3.2) ensures that the iterate avoids the set where the unstable stationary points lie, and hence efficiently contracts to the basin of attraction of the independent component pairs.

Analogous to Theorem 4, when the sample size satisfies the scaling condition  $\frac{1}{d} = \epsilon^{-1} d^3$ , one can carefully choose a stepsize  $\epsilon = \epsilon(d \log T = T)$  and establish a finite-sample bound  $\bar{\epsilon}_T$  based on Lemma 6. We formulate this fact as our second main theorem.

**Theorem 7 (Finite-Sample Error Bound with Uniform Initialization)** Let the dimension  $d \geq 2$ , let Assumption 2 hold, and let initialization  $u^{(0)}$  be uniformly sampled from the unit sphere  $\mathbb{S}^{d-1}$ .

For sample size  $T \geq 100$ , set the stepsize as  $\epsilon(T) = \frac{4d \log \frac{j-4-3j^2}{4B^8 d} T}{j-4-3jT}$ : Then, for any fixed positive numbers  $\delta \geq 2; T \geq 100; \eta \in (0; 1=4]$  satisfying the scaling condition

$$d \geq 2^{\frac{p}{2}} \frac{1}{e} \log \frac{1}{\eta} + 1; \quad C_{7;T} \log^8(C_{7;T}^0 \frac{1}{d} T) \frac{B^8}{j-4-3j^2} \frac{d^3 \log^2 d \log^2 T}{T} \leq \frac{\eta^2}{\log^2 \frac{1}{\eta}}; \quad (3.4)$$

there exists a uniformly distributed random variable  $\epsilon \in [0; 1]$  and an event  $H_{7;T}$  with  $P(H_{7;T}) \geq 1 - \frac{1}{4}$  such that on  $H_{7;T}$ , iteration  $u^{(t)}$  of Algorithm 1 satisfies

$$\tan \setminus u^{(T)}; a_i \leq C_{7;T} \log^{5=2}(C_{7;T}^0 \frac{1}{d}) \frac{B^4}{j-4-3j} \frac{\epsilon^{\frac{p}{2}} d^4 \log^2 T}{T};$$

where  $C_{7;T}; C_{7;T}^0; C_{7;T}^0$  are positive, absolute constants.

Theorem 7 achieves, under the scaling condition  $\frac{1}{d} = \epsilon^{-1} d^3$ , an  $\mathcal{O}(\frac{1}{d^4} \sqrt{\frac{p}{T}})$  finite-sample error bound on  $\tan \setminus u^{(T)}; a_i$  for some  $i$  drawn uniformly at random in  $[d]$ . With Theorems 4 and 7 for warm and uniform initializations respectively, a specific choice of stepsize allows us to have the best of the two worlds. Assuming prior knowledge of the sample size, we initialize  $u^{(0)}$  uniformly at random from the unit sphere  $\mathbb{S}^{d-1}$  and run Algorithm 1 in two consecutive phases, each using  $T=2$  observations:

In the first phase, we initialize  $u^{(0)}$  uniformly at random on unit sphere  $\mathbb{S}^{d-1}$ , pick a constant stepsize  $\epsilon_1 = \epsilon(d \log T = T)$  and update iteration  $u^{(t)}$  via (1.3) for  $T=2$  iterates. Theorem 7 guarantees with high probability that  $u^{(T=2)}$  satisfies the warm initialization condition (2.1) under the scaling condition  $\frac{1}{d} = \epsilon^{-1} d^4$ ;

In the second phase, we warm-initialize the algorithm using the output of the first phase  $u^{(T=2)}$ , pick a constant stepsize  $\epsilon_2 = \epsilon(\log T = T)$  and update the iteration  $u^{(t)}$  via (1.3) for  $T=2$  iterates. The last iterate achieves an error bound  $\bar{\epsilon}_T$  (of  $d = T$ ) as indicated by Theorem 4.

This two-phase procedure yields an improved finite-sample error bound  $\bar{\epsilon}_T$  (of  $d = T$ ) under the uniform initialization and scaling condition  $\frac{1}{d} = \epsilon^{-1} d^4$ , formally stated in the following corollary.

**Corollary 8 (Improved Finite-Sample Error Bound with Uniform Initialization)** Let the dimension  $d \geq 2$ , let Assumption 2 hold, and let initialization  $u^{(0)}$  be uniformly sampled from the unit



sphere  $\mathbb{D}_1$ . Set for sample size  $\bar{T} = 200$  the stepsizes as

$$\eta_1(T) = \frac{8d \log \frac{j^4 3j^2 T}{8B^8 d}}{j^4 3jT}; \quad \eta_2(T) = \frac{9 \log \frac{j^4 3j^2 T}{9B^8}}{j^4 3jT}; \quad (3.5)$$

Then, for any fixed positive  $\epsilon \in (0, 1/5]$  satisfying the scaling condition

$$d \leq \frac{\rho}{2} \frac{1}{2e} \log^{-1} \frac{1}{\epsilon}; \quad \epsilon \geq \max\{C_{4,T}; C_{7,T}; C_{7,T}^2\} \log^8(C_{7,T}^0) \frac{B^8}{j^4 3j^2} \frac{d^4 \log^2 T}{T} \frac{1}{\log^2 \frac{1}{\epsilon}}; \quad (3.6)$$

there exists a uniformly distributed random variable  $\mathbf{u}^{(T)} \in \mathbb{D}_1$  and an event  $\mathcal{H}_{8,C}$  with  $P(\mathcal{H}_{8,C}) \geq 1 - \epsilon$  such that on the event  $\mathcal{H}_{8,C}$ , running Algorithm 1 for  $T=2$  iterates with stepsize  $\eta_1(T)$  followed by  $T=2$  iterates with stepsize  $\eta_2(T)$  outputs  $\mathbf{u}^{(T)}$  satisfying

$$\|\mathbf{u}^{(T)} - \mathbf{a}_i\|_2 \leq \frac{\rho}{2} C_{4,T} \log^{5/2}(C_{4,T}^0) \frac{B^4}{j^4 3j} \frac{d \log^2 T}{T};$$

where  $C_{4,T}; C_{4,T}^0; C_{7,T}; C_{7,T}^0$  are positive, absolute constants defined earlier in Theorems 4 and 7.

We make several remarks:

- (i) Our dynamics-based analysis of uniformly-initialized online tensorial ICA is different from that of PCA in Li et al. (2018) (see). In particular, we provide a coordinate-pair analysis of the algorithm to cope with the landscape of ICA. A two-phase procedure is, however, required to achieve a rate that is sharp in dimension.
- (ii) Ge et al. (2015) make use of artificial noise injection in their noisy projected SGD algorithm, which includes online tensorial independent component analysis as an application scenario. The analysis provided in the appendix of Ge et al. (2015), however, is worsened by its unrealistic isotropic covariance assumption for the incurred stochastic noise, which is not satisfied by the tensorial ICA algorithm, noise-injected or not. A straightforward extension to the generic noise satisfied by the tensorial ICA condition was claimed, but the analysis of such a case is not available in Ge et al. (2015).
- (iii) By choosing  $\epsilon$  in rate analysis above as a small positive constant,  $\epsilon \leq \frac{1}{2}$ , we achieve a scaling condition such that  $d^4 = T$  is sufficiently small up to a polylogarithmic factor, we achieve the finite-sample convergence rate with probability no less than  $1 - \epsilon$ . In comparison, Ge et al. (2015) achieve  $d = T^{1/4}$  error bound for some  $\epsilon \geq 2$ , under a more stringent assumption that all independent components are almost surely bounded, which indicates a scaling condition of  $d^8 = T \ll 0$  at best. The scaling condition imposed in Theorem 1, however, leads to an error bound that depends polynomially on the inverse unsuccessful probability.

In the rest of this section we study the uniform initialization case and prove Theorem 7. The key idea behind our analysis is that the uniform initialization is sufficiently far from the set where the

4. Our scaling condition is the best known when ignoring the polylogarithmic factors. Optimistically, if analysis in the  $d^4$  or  $d^{3.5}$  state-of-the-art complexity results in Jin et al. (2021); Fang et al. (2019) can be carried over to Riemannian optimization, setting  $\rho = 1 = d$  implies a complexity of  $d^9$  or  $d^8$  at best (for the ICA objective the smallest Hessian eigenvalue changes from  $1/d$  to  $(1)$  so the Hessian-Lipschitz constant  $(1/d)$  and the stochastic gradient admits a variance  $\mathcal{O}(d)$ ).

unstable stationary points lie (with high probability), and thus a delicate concentration analysis can show that saddle-point avoidance is guaranteed throughout the entire online tensorial ICA algorithm. Inherited from the warm-initialization analysis [2], we recall the rotated iteration  $v^{(t)} = PA^> u^{(t)}$  and the rotated observations  $Y^{(t)} = PA^> X^{(t)}$ , so our online tensorial ICA update rule can still be translated into (2.8). Here the permutation matrix is  $P(I; 1) = P(1; I) = 1 = P(j; j)$  for  $j \in \{1, \dots, d\}$  and 0 elsewhere, in which  $\arg\min_{i \in [2; d]} \tan^2 \angle u^{(0)}; a_i$ . We let the coordinate-wise intermediate initialization region for each  $k \in [2; d]$  be

$$D_{\text{mid};k} = \{v \in \mathbb{R}^d : v_1^2 \leq \max_{i \in [2; d]} v_i^2 \text{ and } v_1^2 \geq 3v_k^2\} \quad (3.7)$$

In addition, we let the cold initialization region be

$$D_{\text{cold}} = \{v \in \mathbb{R}^d : v_1^2 \leq \max_{i \in [2; d]} v_i^2\} \quad (3.8)$$

By definition of the rotated iteration  $v^{(t)}$  and index  $k$ , we know that  $v^{(0)} \in D_{\text{cold}}$  always holds. As we will see later in Lemmas 9, 10 and 11 of this section, when initialized uniformly at random on the unit sphere  $\mathbb{S}^d$  a gap between  $(v_1^{(0)})^2$  and  $\max_{k \in [2; d]} (v_k^{(0)})^2$  persists on a high-probability event  $(H_{10;L})$  as the data dimension  $n \rightarrow \infty$ . Moreover on a high-probability event  $(H_{k \in [2; d]; 9; L})$ , the iterate  $v^{(t)}$  enters the intersection of intermediate regions  $\bigcap_{k \in [2; d]} D_{\text{mid};k}$  within  $T_{0;0.5}^o$  iterations. After  $v^{(t)}$  enters  $\bigcap_{k \in [2; d]} D_{\text{mid};k}$ , on a third high-probability event  $(H_{k \in [2; d]; 11; L})$  it decays exponentially fast and stays within  $\Theta(\frac{1}{\sqrt{d^3}})$ -neighborhood of the independent component pair  $e_1$ .

### 3.1. Initialization in the Intermediate Region

Recall the intermediate initialization region  $D_{\text{mid};k}$  defined in (3.7) for each  $k \in [2; d]$ . We also define a slightly larger coordinate-wise intermediate auxiliary region for each  $k \in [2; d]$  as  $D_{\text{mid-aux};k} = \{v \in \mathbb{R}^d : v_1^2 \leq \max_{i \in [2; d]} v_i^2 \text{ and } v_1^2 \geq 2v_k^2\}$ . When  $v^{(0)} \in D_{\text{mid};k}$ , we define the first time the iterate exits  $D_{\text{mid-aux};k}$  as

$$T_{w;k} = \inf \{t \geq 1 : v^{(t)} \notin D_{\text{mid-aux};k}^c\} \quad (3.9)$$

Thus, for each  $k \in [2; d]$ ,  $T_{w;k}$  is a stopping time with respect to  $\text{Itratio}_t = (v^{(0)}; Y^{(1)}; \dots; Y^{(t)})$  (we suppose that all that appear in the rest of our discussion satisfy  $k \in [2; d]$ , unless stated otherwise).

Our goal is to prove the following high-probability bound for each coordinate  $k$ . For the analysis of the intermediate iterates, we recall the notation  $u_k^{(t)} = v_k^{(t)} = v_1^{(t)}$  previously defined in (2.11).

**Lemma 9** Let the settings in Lemma 6 hold, and  $x$  the coordinate  $k \in [2; d]$  and value  $\epsilon > 0$ . Then for any positive numbers  $\delta$  satisfying the scaling condition (3.2), and given the coordinate-wise intermediate initialization condition  $v^{(0)} \in D_{\text{mid};k}$ , there exists an event  $(H_{k;9;L})$  satisfying  $P(H_{k;9;L}) \geq 1 - \epsilon - 12 + \frac{5184}{\log^5 \frac{1}{\delta}}$ ; such that on event  $(H_{k;9;L})$  the following holds:

$$\sup_{t \in T_{0;0.5}^o \wedge T_{w;k}} U_k^{(t)} \leq U_k^{(0)} \prod_{s=0}^{t-1} \frac{1}{1 - \frac{1}{3j} (v_1^{(s)})^2 - (v_k^{(s)})^2} \leq 2C_{9;L} \log^{5=2} \frac{1}{\delta} B^4 d^{1=2} (T_{0;0.5}^o)^{1=2}; \quad (3.10)$$

where  $C_{9;L}$  is a positive, absolute constant.

From (3.10) in Lemma 9 we know that, for each coordinate  $k \in [2; d]$  initialized in the intermediate region  $D_{\text{mid};k}$ , with high probability the iteration  $W_k^{(t)}$  fluctuates around a deterministic curve that decays geometrically whenever  $(v_1^{(t)})^2 - (v_k^{(t)})^2$  is bounded below by a positive number.

### 3.2. Initialization in the Cold Region

Recall the cold initialization region  $D_{\text{cold}}$  defined earlier in (3.8). The iteration we study in the cold initialization analysis is

$$W_k^{(t)} = \frac{(v_1^{(t)})^2 - (v_k^{(t)})^2}{(v_k^{(t)})^2}. \quad (3.11)$$

Under the setting of uniform initialization on the unit sphere in Lemma 6, we have the following lemma. Note that the uniform initialization conditions for  $(\rho)$  and  $v^{(0)}$  are equivalent.

Lemma 10 Let  $v^{(0)}$  be uniformly sampled from the unit sphere and  $\epsilon$  be any fixed positive number, with dimension  $d$  and  $\rho$  satisfying

$$d \geq \frac{\rho}{2\epsilon} \log \frac{1}{1-\epsilon}. \quad (3.12)$$

Then there exists an event  $H_{10;L}$  with  $P(H_{10;L}) \geq 1 - 3\epsilon$  such that on event  $H_{10;L}$  the following holds:

$$\min_{k \in [2; d]} W_k^{(0)} \geq \frac{\epsilon}{8 \log \frac{1}{1-\epsilon} \log d}. \quad (3.13)$$

Our goal is to estimate the time when  $v^{(t)}$  enters each coordinate-wise intermediate initialization region starting with the initialization gap given in Lemma 10. For each coordinate  $k \in [2; d]$ , we define the first time  $v^{(t)}$  enters the coordinate-wise intermediate region  $D_{\text{mid};k}$  as

$$T_{C;k} = \inf_{t \geq 1} \{v^{(t)} \in D_{\text{mid};k}\} \quad (3.14)$$

and the first time the iterate exits  $D_{\text{cold}}$  without entering  $D_{\text{mid};k}$ , earlier defined in (3.7) and (3.8), as

$$T_1 = \inf_{t \geq 1} \{v^{(t)} \notin D_{\text{cold}}^c\} \quad (3.15)$$

In other words,  $T_{C;k}$  is the first time  $t$  such that  $W_k^{(t)} \geq 2$ , and  $T_1$  is the first time  $\min_{i \in [2; d]} W_i^{(t)} < 0$ . By the definition of the rotated iteration  $v^{(t)}$  we have  $v^{(0)} \in D_{\text{cold}}$  always holds. For each coordinate  $k \in [2; d]$ , if  $T_{C;k} = 0$  then  $v^{(0)} \in D_{\text{mid};k}$  and the previous analysis directly applies for coordinate  $k$ . The following lemma characterizes the opposite case, where the exponential growth of iteration  $W_k^{(t)}$  helps us to determine when  $v^{(t)}$  enters the intermediate region  $D_{\text{mid};k}$ .

Lemma 11 Let the settings in Lemma 6 hold, and  $k \in [2; d]$ . Then, for any fixed positive numbers  $\epsilon$  satisfying the scaling condition (3.2) along with the coordinate-wise

cold initialization condition  $v^{(0)} \notin D_{\text{cold}} \setminus D_{\text{mid};k}^c$ , there exists an event  $\mathcal{H}_{k;11;L}$  with  $P(\mathcal{H}_{k;11;L})$

$1 - 15 + \frac{5184}{\log^5 - 1}$ ; such that on event  $\mathcal{H}_{k;11;L}$  the following holds:

$$\sup_{t \in T_{0;0.5}^o \wedge T_{c;k} \wedge T_1} W_k^{(t)} \leq \frac{1}{1 + j^4} \frac{3j(v_1^{(s)})^2}{1} W_k^{(0)} \leq C_{11;L} \log^{5=2} \frac{1}{j^4} \frac{B^4}{3j^{1=2}} d^{1=2}; \quad (3.16)$$

where  $C_{11;L}$  is a positive, absolute constant.

From Lemma 11 we know that for each coordinate  $k \in [2; d]$ , if the initialization lies outside intermediate region  $D_{\text{mid};k}$ , then, with high probability, the iterate  $v_k^{(t)}$  is controlled within an exponentially growing dynamics for  $T_{0;0.5}^o$  iterates before it either enters the intermediate region  $D_{\text{mid};k}$  or exits the cold region  $D_{\text{cold}}$ . As we will see later in the proof of Lemma 6, by putting all coordinates together we can show that  $v^{(t)}$  rarely leaves the cold region  $D_{\text{cold}}$ , which implies that  $v^{(t)}$  will enter the joint intermediate region  $\cap_{k \in [2; d]} D_{\text{mid};k}$  within  $T_{0;0.5}^o$  iterates with high probability.

#### 4. Additional Related Literature

The themes of ICA and tensor decomposition have been studied in numerous statistical and signal-processing literatures (Bach and Jordan, 2002; Chen and Bickel, 2006; Samworth and Yuan, 2012; Bonhomme and Robin, 2009; Eriksson and Koivunen, 2004; Hallin and Mehta, 2015; Arihoy et al., 2001; Hyärinen and Oja, 1997; Hyvarinen, 1999; Hyvarinen and Oja, 2000; Ilmonen and Paindaveine, 2011; Kollo, 2008; Miettinen et al., 2015; Oja et al., 2006; Tichavsky et al., 2006; Wang and Lu, 2017; Ge and Ma, 2017). For a treatment from a spectral learning perspective (mainly for the deterministic scenario), we refer to the recently published monograph of Janzamin et al. (2019) and the bibliography therein. Recent literature studies the ICA setting in the context of specific parametric families for independent component distributions and shows that parametric (Lee et al., 1999), semi-parametric (Hastie and Tibshirani, 2003; Chen and Bickel, 2006; Ilmonen and Paindaveine, 2011) and nonparametric (Bach and Jordan, 2002; Samarov and Tsybakov, 2004; Samworth and Yuan, 2012) models can be estimated via maximal likelihood estimation or minimization of mutual information between independent components. Our work focuses on a different type of contrast function based on tensor decomposition and kurtosis maximization, and hence our methodology is quite different from this line of work.

Our work is most closely related to Ge et al. (2015), which led to a general line of work on stochastic-gradient-based nonconvex optimization (Dauphin et al., 2014; Ge et al., 2015; Sun et al., 2015, 2017, 2018; Anandkumar and Ge, 2016; Jin et al., 2017, 2021, 2018; Lei et al., 2017; Allen-Zhu, 2018; Daneshmand et al., 2018; Fang et al., 2019; Cutkosky and Orabona, 2019; Cutkosky and Mehta, 2020) as well as the Riemannian manifold (Zhang and Sra, 2016; Zhang et al., 2016; Tripuraneni et al., 2018). A large family of nonconvex landscape has been studied in Mei et al. (2018), where uniform convergence of the empirical loss to the population loss is established. It is also related to work on recursive variance-reduced gradient methods for smooth optimization (Nguyen et al., 2017a,b; Fang et al., 2018; Zhou et al., 2020; Wang et al., 2019; Arjevani et al., 2020). In particular, we note the work of Ge et al. (2015); Jin et al. (2021); Daneshmand et al. (2018); Fang et al. (2019), who study the dynamics of SGD for optimizing generic functions. When applied to specific statistical estimation problems, however, the results obtained by these methods can be coarse due to their neglect of specific geometric features of the landscape.

The pioneering work of [Ge et al. \(2015\)](#) studied the convergence rate of SGD for minimizing a large class of nonconvex objectives defined on a generic Riemannian manifold. Under the bounded distributional assumption, [Ge et al. \(2015\)](#) prove that SGD equipped with projection as well as a special noise-injection step can escape from all saddle points and land at an approximate local minimizer in polynomial time of relevant parameters. Convergence rates for generic first-order gradient descent algorithms without adding noise injection are generally unknown and can be unfavorable ([Lee et al., 2019](#); [Du et al., 2017](#); [Pemantle, 1990](#)). Favorable results can be obtained under special spherical distributions for the noise or sophisticated procedures for avoidance of saddle points ([Ge et al., 2015](#); [Jin et al., 2017, 2021](#); [Sun et al., 2015](#); [Allen-Zhu, 2018](#); [Fang et al., 2019](#)).

Recent years have witnessed significant progress on computational and statistical aspects of low-rank representation methods. A number of recent papers ([Carmon and Duchi, 2020](#); [Ge et al., 2017](#); [Li et al., 2018](#); [Bai et al., 2018](#); [Davis et al., 2018](#); [Li and Bresler, 2018](#); [Zhu et al., 2018](#); [Ma et al., 2019](#); [Chen et al., 2019](#); [Gilboa et al., 2019](#); [Kuo et al., 2019](#); [Qu et al., 2019](#); [Tan and Vershynin, 2019](#); [Yang et al., 2019](#); [Na et al., 2019](#); [Chen et al., 2020](#); [Lau et al., 2020](#); [Li et al., 2020](#); [Zhai et al., 2020](#)) study the gradient-descent dynamics of matrix factorization/completion, principal component pursuit, dictionary learning, phase retrieval, blind deconvolution, and many others, in the setting of batch or online (streaming) data. Notably, [Carmon and Duchi \(2020\)](#); [Li et al. \(2018\)](#) pursue a dynamics-based analysis to study gradient descent and its cubic-regularized variant for eigenvalue problems. Related work on efficient convergence of Oja's online PCA iteration can be found in [Jain et al. \(2016\)](#); [Allen-Zhu and Li \(2017\)](#); [Li et al. \(2016\)](#); [Wang and Wu \(2020\)](#); [Wang and Lu \(2017\)](#), who study the online tensorial ICA method from the viewpoint of scaling limits and (stochastic) differential equation approximations. While these results provide valuable insights into our problem, straightforward translations to nonasymptotic convergence guarantees are not available, mainly due to the differential equation approximation being a weak convergence formulation instead of a strong one. Our refined projected stochastic gradient analysis for online tensorial ICA provided in this paper is both dynamics-based and nonasymptotic, and we are able to prove that under some mild scaling conditions, random initialization provides sufficient deviation from the set of unstable stationary points such that a vanilla online tensorial ICA algorithm can be guaranteed to achieve a sharp convergence rate.

## 5. Summary

We have studied the dynamics of an algorithm formulated as online stochastic approximation of (orthogonal) tensorial ICA. Our algorithm can be viewed as a method for optimizing a nonconvex objective of excess kurtosis in a given direction. We show that with properly chosen stepsizes and under mild scaling conditions our online tensorial ICA algorithm achieves  $(\epsilon \text{ and } \bar{T})$ -convergence rate, which is superior to the best existing analysis of this problem. Our algorithm requires no noise-injection steps or specially-designed loops for saddle-point avoidance, and our dynamics-based approach enjoys multiple advantages over existing analyses of online stochastic approximation for tensorial ICA estimation.

We believe that our analysis can generalize to a broader class of statistical estimation problems that can be cast as nonconvex stochastic optimization problems. Future directions include further improvements of the convergence rate and scaling conditions or justification of the impossibility (or minimax optimality) of such rates, analyzing the mini-batch stochastic approximation algorithm as well as the non-identical kurtosis case for ICA, and finally generalizing our analysis of

the dynamics of stochastic online algorithms to the nonorthogonal tensor decomposition case and over-parameterized cases.

## Acknowledgments

We thank Wenlong Mou and Yuren Zhou for valuable discussions. This work was supported in part by the Mathematical Data Science program of the Office of Naval Research under grant number N00014-18-1-2764.

## References

- Zeyuan Allen-Zhu. Natasha 2: Faster non-convex optimization than SGD. *Advances in Neural Information Processing Systems*, pages 2676–2687, 2018.
- Zeyuan Allen-Zhu and Yuanzhi Li. First efficient convergence for streaming PCA: a global, gap-free, and near-optimal rate. *The 58th Annual Symposium on Foundations of Computer Science* 2017.
- Animashree Anandkumar and Rong Ge. Efficient approaches for escaping higher order saddle points in non-convex optimization. *Conference on Learning Theory*, pages 81–102, 2016.
- Yossi Arjevani, Yair Carmon, John C Duchi, Dylan J Foster, Ayush Sekhari, and Karthik Sridharan. Second-order information in non-convex stochastic optimization: Power and limitations. In *Conference on Learning Theory*, pages 242–299. PMLR, 2020.
- Francis R Bach and Michael I Jordan. Kernel independent component analysis. *Journal of Machine Learning Research*, 3:1–48, 2002.
- Yu Bai, Qijia Jiang, and Ju Sun. Subgradient descent learns orthogonal dictionary. *International Conference on Learning Representations*, 2018.
- Stéphane Bonhomme and Jean-Marc Robin. Consistent noisy independent component analysis. *Journal of Econometrics*, 149(1):12–25, 2009.
- Yair Carmon and John C Duchi. First-order methods for nonconvex quadratic minimization. *SIAM Review* 62(2):395–436, 2020.
- Aiyou Chen and Peter J Bickel. Efficient independent component analysis. *The Annals of Statistics* 34(6):2825–2855, 2006.
- Yuxin Chen, Yuejie Chi, Jianqing Fan, and Cong Ma. Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval. *Mathematical Programming* 176(1):5–37, 2019.
- Yuxin Chen, Yuejie Chi, Jianqing Fan, and Cong Ma. Spectral methods for data science: A statistical perspective. *arXiv preprint arXiv:2012.08496*, 2020.
- Pierre Comon. Independent component analysis, a new concept. *Signal Processing* 36(3):287–314, 1994.

- Ashok Cutkosky and Harsh Mehta. Momentum improves normalized SGD. *International Conference on Machine Learning*, pages 2260–2268. PMLR, 2020.
- Ashok Cutkosky and Francesco Orabona. Momentum-based variance reduction in non-convex SGD. In *Advances in Neural Information Processing Systems*, pages 15236–15245, 2019.
- Hadi Daneshmand, Jonas Kohler, Aurelien Lucchi, and Thomas Hofmann. Escaping saddles with stochastic gradients. *International Conference on Machine Learning*, pages 1163–1172, 2018.
- Yann N Dauphin, Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, Surya Ganguli, and Yoshua Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. In *Advances in Neural Information Processing Systems*, pages 2933–2941, 2014.
- Damek Davis, Dmitriy Drusvyatskiy, Kellie J MacPhee, and Courtney Paquette. Subgradient methods for sharp weakly convex functions. *Journal of Optimization Theory and Applications*, 179(3):962–982, 2018.
- Simon S Du, Chi Jin, Jason D Lee, Michael I Jordan, Bálint Póczos, and Aarti Singh. Gradient descent can take exponential time to escape saddle points. *Advances in Neural Information Processing Systems*, pages 1068–1078, 2017.
- Jan Eriksson and Visa Koivunen. Identifiability, separability, and uniqueness of linear ICA models. *IEEE Signal Processing Letters*, 1(7):601–604, 2004.
- Xiequan Fan, Ion Grama, and Quansheng Liu. Large deviation exponential inequalities for supermartingales. *Electronic Communications in Probability*, 17, 2012.
- Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. SPIDER: Near-optimal non-convex optimization via stochastic path integrated differential estimation. *Advances in Neural Information Processing Systems*, pages 686–696, 2018.
- Cong Fang, Zhouchen Lin, and Tong Zhang. Sharp analysis for nonconvex SGD escaping from saddle points. In *Conference on Learning Theory*, pages 1192–1234. PMLR, 2019.
- Alan Frieze, Mark Jerrum, and Ravi Kannan. Learning linear transformations. *Proceedings of 37th Conference on Foundations of Computer Science*, pages 359–368. IEEE, 1996.
- Rong Ge and Tengyu Ma. On the optimization landscape of tensor decomposition. *Advances in Neural Information Processing Systems*, pages 3653–3663, 2017.
- Rong Ge, Furong Huang, Chi Jin, and Yang Yuan. Escaping from saddle points – online stochastic gradient for tensor decomposition. *Conference on Learning Theory*, pages 797–842, 2015.
- Rong Ge, Chi Jin, and Yi Zheng. No spurious local minima in nonconvex low rank problems: A unified geometric analysis. *International Conference on Machine Learning*, pages 1233–1242. PMLR, 2017.
- Dar Gilboa, Sam Buchanan, and John Wright. Efficient dictionary learning with gradient descent. In *International Conference on Machine Learning*, pages 2252–2259. PMLR, 2019.

- Marc Hallin and Chintan Mehta. R-estimation for asymmetric independent component analysis. *Journal of the American Statistical Association* 110(509):218–232, 2015.
- Trevor Hastie and Rob Tibshirani. Independent components analysis through product density estimation. In *Advances in Neural Information Processing Systems*, pages 665–672, 2003.
- Aapo Hyvärinen. Fast and robust xed-point algorithms for independent component analysis. *Transactions on Neural Networks* 10(3):626–634, 1999.
- Aapo Hyvärinen and Erkki Oja. A fast xed-point algorithm for independent component analysis. *Neural Computation* 9(7):1483–1492, 1997.
- Aapo Hyvärinen and Erkki Oja. Independent component analysis: algorithms and applications. *Neural Networks* 13(4-5):411–430, 2000.
- Aapo Hyvärinen, Juha Karhunen, and Erkki Oja. *Independent Component Analysis*. John Wiley & Sons, 2001.
- Pauliina Ilmonen and Davy Paindaveine. Semiparametrically efficient inference based on signed ranks in symmetric independent component models. *The Annals of Statistics* 39(5):2448–2476, 2011.
- Prateek Jain, Chi Jin, Sham M Kakade, Praneeth Netrapalli, and Aaron Sidford. Streaming PCA: Matching matrix Bernstein and near-optimal finite sample guarantees for Oja's algorithm. In *Conference on Learning Theory*, pages 1147–1164. PMLR, 2016.
- Majid Janzamin, Rong Ge, Jean Kossai, and Anima Anandkumar. Spectral learning on matrices and tensors. *Foundations and Trends in Machine Learning* 12(5-6):393–536, 2019.
- Chi Jin, Rong Ge, Praneeth Netrapalli, Sham M Kakade, and Michael I Jordan. How to escape saddle points efficiently. In *International Conference on Machine Learning*, pages 1724–1732. PMLR, 2017.
- Chi Jin, Praneeth Netrapalli, and Michael I Jordan. Accelerated gradient descent escapes saddle points faster than gradient descent. *Conference on Learning Theory*, pages 1042–1085, 2018.
- Chi Jin, Praneeth Netrapalli, Rong Ge, Sham M Kakade, and Michael I Jordan. On nonconvex optimization for machine learning: Gradients, stochasticity, and saddle points. *Journal of the ACM*, 68(2):1–29, 2021.
- Tõnu Kollo. Multivariate skewness and kurtosis measures with an application in JCA. *Journal of Multivariate Analysis* 99(10):2328–2338, 2008.
- Han-Wen Kuo, Yenson Lau, Yuqian Zhang, and John Wright. Geometry and symmetry in short-and-sparse deconvolution. In *International Conference on Machine Learning*, pages 3570–3580. PMLR, 2019.
- Nâmane Laib. Exponential-type inequalities for martingale difference sequences. application to nonparametric regression estimation. *Communications in Statistics-Theory and Methods* 28(7):1565–1576, 1999.



- Yenson Lau, Qing Qu, Han-Wen Kuo, Pengchang Zhou, Yuqian Zhang, and John Wright. Short and sparse deconvolution—a geometric approach. *International Conference on Learning Representations*, 2020.
- Jason D Lee, Ioannis Panageas, Georgios Piliouras, Max Simchowitz, Michael I Jordan, and Benjamin Recht. First-order methods almost always avoid strict saddle points. *Mathematical Programming* 176(1):311–337, 2019.
- Te-Won Lee, Mark Girolami, and Terrence J Sejnowski. Independent component analysis using an extended infomax algorithm for mixed subgaussian and supergaussian sources. *Neural Computation*, 11(2):417–441, 1999.
- Lihua Lei, Cheng Ju, Jianbo Chen, and Michael I Jordan. Non-convex nite-sum optimization via SCSSG methods. *Advances in Neural Information Processing Systems*, pages 2348–2358, 2017.
- Emmanuel Lesigne and Dalibor Volyn. Large deviations for martingales. *Stochastic Processes and Their Applications* 96(1):143–159, 2001.
- Chris Junchi Li, Zhaoran Wang, and Han Liu. Online ICA: Understanding global dynamics of nonconvex optimization via diffusion processes. *Advances in Neural Information Processing Systems*, pages 4967–4975, 2016.
- Chris Junchi Li, Mengdi Wang, Han Liu, and Tong Zhang. Near-optimal stochastic approximation for online principal component estimation. *Mathematical Programming* 167(1):75–97, 2018.
- Xiao Li, Zhihui Zhu, Anthony Man-Cho So, and Rene Vidal. Nonconvex robust low-rank matrix recovery. *SIAM Journal on Optimization* 30(1):660–686, 2020.
- Yanjun Li and Yoram Bresler. Global geometry of multichannel sparse blind deconvolution on the sphere. *Advances in Neural Information Processing Systems*, pages 1132–1143, 2018.
- Cong Ma, Kaizheng Wang, Yuejie Chi, and Yuxin Chen. Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion, and blind deconvolution. *Foundations of Computational Mathematics* 20:19, 2019.
- Song Mei, Yu Bai, and Andrea Montanari. The landscape of empirical risk for nonconvex losses. *The Annals of Statistics* 46(6A):2747–2774, 2018.
- Jari Miettinen, Sara Taskinen, Klaus Nordhausen, and Hannu Oja. Fourth moments and independent component analysis. *Statistical Science* 30(3):372–390, 2015.
- Sen Na, Zhuoran Yang, Zhaoran Wang, and Mladen Kolar. High-dimensional varying index coefficient models via Stein's identity. *Journal of Machine Learning Research* 20:1–44, 2019.
- Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Tak. SARAH: A novel method for machine learning problems using stochastic recursive gradient. *International Conference on Machine Learning*, pages 2613–2621, 2017a.
- Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Tak. Stochastic recursive gradient algorithm for nonconvex optimization. *arXiv preprint arXiv:1705.07261*, 2017b.

- Hannu Oja, Seija Sirin, and Jan Eriksson. Scatter matrices and independent component analysis. *Austrian Journal of Statistics* 35(2&3):175–189, 2006.
- Robin Pemantle. Nonconvergence to unstable points in urn models and stochastic approximations. *The Annals of Probability* pages 698–712, 1990.
- Qing Qu, Xiao Li, and Zihui Zhu. A nonconvex approach for exact and efficient multichannel sparse blind deconvolution. *Advances in Neural Information Processing Systems* pages 4015–4026, 2019.
- Alexander Samarov and Alexandre Tsybakov. Nonparametric independent component analysis. *Bernoulli*, 10(4):565–582, 2004.
- Richard J Samworth and Ming Yuan. Independent component analysis via nonparametric maximum likelihood estimation. *The Annals of Statistics* 40(6):2973–3002, 2012.
- James V Stone. *Independent Component Analysis: A Tutorial Introduction*. MIT press, 2004.
- Ju Sun, Qing Qu, and John Wright. When are nonconvex problems not so hard? preprint arXiv:1510.06096 2015.
- Ju Sun, Qing Qu, and John Wright. Complete dictionary recovery over the sphere I: Overview and the geometric picture. *IEEE Transactions on Information Theory* 63(2):853–884, 2017.
- Ju Sun, Qing Qu, and John Wright. A geometric analysis of phase retrieval. *Foundations of Computational Mathematics* 18(5):1131–1198, 2018.
- Yan Shuo Tan and Roman Vershynin. Online stochastic gradient descent with arbitrary initialization solves non-smooth, non-convex phase retrieval. preprint arXiv:1910.12837 2019.
- Petr Tichavsky, Zbynek Koldovsky, and Erkki Oja. Performance analysis of the FastICA algorithm and crame/r-rao bounds for linear independent component analysis. *IEEE Transactions on Signal Processing* 54(4):1189–1203, 2006.
- Nilesh Tripuraneni, Mitchell Stern, Chi Jin, Jeffrey Regier, and Michael I Jordan. Stochastic cubic regularization for fast nonconvex optimization. *Advances in Neural Information Processing Systems* pages 2899–2908, 2018.
- Aad van der Vaart and Jon A Wellner. *Weak Convergence and Empirical Processes: With Application to Statistics*. Springer, 1996.
- Vincent Q Vu and Jing Lei. Minimax sparse principal subspace estimation in high dimensions. *Annals of Statistics* 41(6):2905–2947, 2013.
- Martin J Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*, volume 48. Cambridge University Press, 2019.
- Chuang Wang and Yue Lu. The scaling limit of high-dimensional online independent component analysis. *Advances in Neural Information Processing Systems* pages 6638–6647, 2017.

- Yazhen Wang and Shang Wu. Asymptotic analysis via stochastic differential equations of gradient descent algorithms in statistical and computational paradigms. *Journal of Machine Learning Research* 21(199):1–103, 2020.
- Zhe Wang, Kaiyi Ji, Yi Zhou, Yingbin Liang, and Vahid Tarokh. SpiderBoost and momentum: Faster variance reduction algorithms. *Advances in Neural Information Processing Systems* pages 2406–2416, 2019.
- Zhuoran Yang, Lin F Yang, Ethan X Fang, Tuo Zhao, Zhaoran Wang, and Matey Neykov. Misspecified nonconvex statistical optimization for sparse phase retrieval. *Mathematical Programming* 176(1):545–571, 2019.
- Yuexiang Zhai, Zitong Yang, Zhenyu Liao, John Wright, and Yi Ma. Complete dictionary learning via  $\ell^4$ -norm maximization over the orthogonal group. *Journal of Machine Learning Research* 21(165):1–68, 2020.
- Hongyi Zhang and Suvrit Sra. First-order methods for geodesically convex optimization. *Conference on Learning Theory* pages 1617–1638, 2016.
- Hongyi Zhang, Sashank J Reddi, and Suvrit Sra. Riemannian SVRG: Fast stochastic optimization on Riemannian manifolds. *Advances in Neural Information Processing Systems* pages 4592–4600, 2016.
- Dongruo Zhou, Pan Xu, and Quanquan Gu. Stochastic nested variance reduction for nonconvex optimization. *Journal of Machine Learning Research* 21(103):1–63, 2020.
- Zhihui Zhu, Yifan Wang, Daniel Robinson, Daniel Naiman, Rene Vidal, and Manolis C Tsakiris. Dual principal component pursuit: Improved analysis and efficient algorithms. *Advances in Neural Information Processing Systems* pages 2175–2185, 2018.

## Appendix

In Appendix, [x](#) [A](#) proves the main results in the paper, [x](#) [B](#) and [x](#) [C](#) provide all secondary lemmas and their proofs for warm and uniform initialization analysis, separately. [x](#) [E](#) and [x](#) [F](#) provide necessary tools including a reversed version of Gronwall's inequality, preliminaries and properties of Orlicz  $\psi$ -norm and a concentration inequality, all of which are theoretical building blocks of this paper. Finally [x](#) [G](#) visualizes the tensorial ICA landscape and provides preliminary simulation results that validate our theory.

**Notations** Throughout this paper, we treat  $\epsilon, \delta, \eta$  as positive constants. We use bold upper case letters to denote matrices, bold lower case letters to denote vectors and italic letters to denote randomness. For any matrix  $A$  or vector  $v$ ,  $A^T$  and  $v^T$  denote their transposes. For any vector  $v_k$  denotes its  $k$ th coordinate. For a sequence of  $x^{(t)}$  and positive  $y^{(t)}$ , we write  $x^{(t)} = O(y^{(t)})$  if there exists a positive constant  $M$  such that  $|x^{(t)}| \leq M y^{(t)}$ , write  $x^{(t)} = \Theta(y^{(t)})$  if there exists a positive constant  $M < 1$  such that  $|x^{(t)}| \leq M y^{(t)}$ , and write  $x^{(t)} = \asymp(y^{(t)})$  if both  $x^{(t)} = O(y^{(t)})$  and  $x^{(t)} = \Omega(y^{(t)})$  hold. We use  $\tilde{O}, \tilde{\Theta}, \tilde{\asymp}$  to hide factors that are polylogarithmically dependent on dimension  $d$ , stepsize  $\eta$ , sample size  $n$  and inverse unsuccessful probability  $1 - \epsilon$ . We use  $\lfloor x \rfloor$  to denote the floor function and  $\lceil x \rceil$  to denote the ceiling function. We let  $x \wedge y = \min(x, y)$  and  $x \vee y = \max(x, y)$ . For any vector  $v$ , we use  $\|v\|$  to denote its Euclidean norm. For any integer we define  $\text{set}[n] = \{1, \dots, n\}$ . Finally, we use  $S^c$  to denote the complement of set (or event).

## Appendix A. Deferred Proofs of Main Results

This section includes the deferred proofs of the main results [x](#) [2](#) and [x](#) [3](#). In their order of appearance [x](#) [A.1](#) and [x](#) [A.2](#) prove in sequel Lemma [3](#) and Theorem [4](#), [x](#) [A.3](#) and [x](#) [A.4](#) prove the convergence Lemma [6](#) and its finite-sample error Theorem [7](#), respectively. All secondary lemmas are deferred to [x](#) [B](#) and [x](#) [C](#).

### A.1. Proof of Lemma [3](#)

We use the key Lemma [5](#) to tightly estimate the dynamics in all coordinates and thereby obtain a proof of Lemma [3](#).

**Proof [Proof of Lemma [3](#)]** We denote  $\mathcal{E}_{3;L} = \bigcap_{k \in [2;d]} \mathcal{E}_{k;5;L}$  as the intersection of events  $\mathcal{E}_{k;5;L}$ . Consider now the following  $(d-1)$ -dimensional vector

$$U_k^{(t)} = U_k^{(0)} + \sum_{s=0}^{t-1} \eta \sum_{j=4}^3 \sum_{i=1}^d (v_1^{(s)})^2 (v_k^{(s)})^2 \mathbb{1}_{\{k \in [2;d]\}} \quad (A.1)$$

Using Lemma [5](#), on the event  $\mathcal{E}_{3;L} \setminus (\mathcal{E}_{3;L}^c \cup \mathcal{T}_x)$  we bound the Euclidean norm of (A.1) by

$$\sqrt{\sum_{k=2}^d \left( U_k^{(t)} - U_k^{(0)} \right)^2} \leq \sum_{k=2}^d \sum_{s=0}^{t-1} \eta \sum_{j=4}^3 \sum_{i=1}^d (v_1^{(s)})^2 (v_k^{(s)})^2 \leq \sqrt{d} \cdot 2C_{5;L} \log^{5=2} \eta B^4 (\mathcal{T}_x)^{1=2} \quad (A.2)$$

Additionally, the left hand of (A.2) is the norm of the subtraction of two vectors and hence lower bounded by

$$\sqrt{\sum_{k=2}^d \|X_k^{(t)} - U_k^{(0)}\|^2} \geq \sqrt{\sum_{k=2}^d \|X_k^{(t)}\|^2 - \sum_{k=2}^d \|U_k^{(0)}\|^2} \quad (A.3)$$

due to triangle inequality of Euclidean norm. The definition of iterative  $U_k^{(t)}$  in (2.11) implies that  $\tan \angle v^{(t)}; e_1 = \frac{\sum_{k=2}^d \|U_k^{(t)}\|^2}{\sum_{k=2}^d \|X_k^{(t)}\|^2}$  and the definition of stopping time  $T_x$  in (2.12) implies that  $(v_1^{(s)})^2 - (v_k^{(s)})^2 \geq 1/3$  holds for all  $k \in [2; d]$  and  $0 \leq s < t$  on the event  $H_{3;L} \cap (t \leq T_x)$ . Combining this with (A.2) and (A.3), we obtain on the event  $H_{3;L} \cap (t \leq T_x)$

$$\begin{aligned} \tan \angle v^{(t)}; e_1 &= \frac{\sum_{k=2}^d \|X_k^{(t)}\|^2}{\sum_{k=2}^d \|U_k^{(t)}\|^2} \\ &\geq \frac{\sum_{k=2}^d \|X_k^{(0)}\|^2 - \sum_{k=2}^d \|U_k^{(0)}\|^2}{\sum_{k=2}^d \|U_k^{(t)}\|^2} \\ &= \tan \angle v^{(0)}; e_1 - \frac{\sum_{k=2}^d \|U_k^{(0)}\|^2}{\sum_{k=2}^d \|U_k^{(t)}\|^2} \end{aligned} \quad (A.4)$$

The definition of  $T_x$  in (2.2) along with  $\log \frac{1}{3} \leq \frac{3j}{j^4 - 3j} \leq \frac{3}{j^4 - 3j}$  gives

$$T_x \leq 1 + \frac{3}{j^4 - 3j} \log \frac{j^4 - 3j}{B^8} \leq \frac{3(+1)}{j^4 - 3j} \log \frac{j^4 - 3j}{B^8} \quad (A.5)$$

where in the second inequality we use  $\frac{j^4 - 3j}{B^8} \leq e$  and  $\frac{1}{j^4 - 3j} \leq 1$  implied by scaling condition (2.3). Using relation (A.5), we find that scaling condition (2.3) with constant  $713C_{5,L}^2$  indicates

$$\begin{aligned} 2C_{5,L} \log^{5-2} \frac{1}{B^4} (d^2 T_x)^{1-2} &\leq \frac{4C_{5,L}^2 B^8 \log^5 \frac{1}{B^4} \frac{3(+1)}{j^4 - 3j} d \log \frac{j^4 - 3j}{B^8}}{713} \\ &\leq \frac{12}{713} (+1) C_{3,L} \log^8(T_{x,1}) \frac{B^8}{j^4 - 3j} d \log \frac{j^4 - 3j}{B^8} \leq \frac{12}{713} \end{aligned} \quad (A.6)$$

where the elementary inequality  $\log^5 \frac{1}{B^4} \log^8(T_{x,1})$  is applied due to  $T_{x,1} \leq 1$  and  $e \leq 1$ . Viewing (2.1) (equivalent to  $\alpha^{(0)} \geq D_{\text{warm}}$  in (2.10)) and (A.4), we have for each  $t$  on the event  $H_{3;L} \cap (T_x \leq T_x) \cap (t \leq T_x)$  that

$$\tan \angle v^{(t)}; e_1 \geq \tan \angle v^{(0)}; e_1 - \frac{12}{713} < \frac{1}{3} + \frac{1}{2} - \frac{1}{3} = \frac{1}{2}; \quad (A.7)$$

where in the first inequality we again apply (2.3). This further indicates that on event  $H_{3;L} \setminus (T_x > T_*)$ , inequality  $\tan^{-1} v^{(T_x)}; e_1 < 1 = \frac{1}{2}$  holds, which contradicts with the fact that  $v^{(T_x)} \geq D_{\text{warm-aux}}^c$  on the same event. Therefore, we have  $H_{3;L} \setminus (T_x > T_*) = \emptyset$ , and equivalently

$$H_{3;L} = H_{3;L} \setminus (T_x > T_*); \tag{A.8}$$

This implies that for all  $t \in [0; T_*]$ , (A.4) holds on the event  $H_{3;L}$ . Plugging in the inequality involving  $T_*$  in (A.5), we have on the event  $H_{3;L}$  that for all  $t \in [0; T_*]$

$$\begin{aligned} \tan^{-1} v^{(t)}; e_1 &\leq \tan^{-1} v^{(0)}; e_1 + 1 - \frac{1}{3} j^4 - 3j^{-1} + 2C_{5;L} B^4 \log^{5-2} \frac{1}{d} (d^{-2} T_*)^{1-2} \\ &\leq \tan^{-1} v^{(0)}; e_1 + 1 - \frac{1}{3} j^4 - 3j^{-1} + \frac{B^4}{j^4 - 3j^{-1}} \frac{1}{d \log \frac{j^4 - 3j^{-1}}{B^8}}; \end{aligned} \tag{A.9}$$

Letting the constant  $C_{3;L} = 2^p \bar{C}_{5;L}$ , the scaling relation (2.9) and the above derivation (A.9) prove that (2.4) holds for all  $t \in [0; T_*]$  on the event  $H_{3;L}$ .

The only left is to estimate the probability of  $H_{3;L}$ . Lemma 5 gives for each  $k \in [2; d]$  the probability of event  $P(H_{k;5;L}) \leq 1 - 6 + 12 + \frac{5184}{\log^5 \frac{1}{d}}$ , and hence elementary union bound calculation gives

$$P(H_{3;L}) = P\left(\bigcap_{k \in [2; d]} H_{k;5;L}\right) \leq 1 - 6 + 12 + \frac{5184}{\log^5 \frac{1}{d}}; \tag{A.10}$$

completing the whole proof of Lemma 3. ■

### A.2. Proof of Theorem 4

Now we turn to the proof of the finite-sample error Theorem 4. The idea is to apply Lemma 3 with appropriate stepsize (as in (2.5)) as well as an appropriate  $\epsilon$  to obtain the finite-sample error bound.

Proof [Proof of Theorem 4]

(1) We first provide an upper bound on  $\bar{\sigma}_T(T)$ . Under scaling condition (2.6) with constant  $C_{4;T} = \max\{90C_{3;L}; 10\}$  and some constant  $\bar{C}_{4;T} > 1$  to be determined later, we have  $\frac{2j^4 - 3j^2}{9B^8} T \leq \epsilon$ . Plugging in  $\epsilon = \bar{\sigma}_T(T)$  from (2.5) to relation (A.5), we have

$$\begin{aligned} \bar{\sigma}_T(T) &\leq \frac{3(\epsilon + 1)}{j^4 - 3j^{-1}} \bar{\sigma}_T(T) + \log \frac{j^4 - 3j^{-1}}{B^8} \bar{\sigma}_T(T) + \frac{1}{9 \log \frac{2j^4 - 3j^2}{9B^8} T} \log \frac{j^4 - 3j^{-1}}{B^8} \bar{\sigma}_T(T) + \frac{2(\epsilon + 1)}{3} T; \\ &= \frac{3(\epsilon + 1)}{j^4 - 3j^{-1}} \frac{2j^4 - 3j^2}{9 \log \frac{2j^4 - 3j^2}{9B^8} T} \bar{\sigma}_T(T) + \log \frac{j^4 - 3j^{-1}}{B^8} \bar{\sigma}_T(T) + \frac{2(\epsilon + 1)}{3} T; \end{aligned} \tag{A.11}$$

(2) Next we provide a lower bound of  $\sigma_{(T)}$ . By Taylor expansion, for all  $x \in (0, 1/3]$  we know that

$$\log(1-x) + x = \sum_{n=2}^{\infty} \frac{x^n}{n} = \frac{x^2}{2} \sum_{n=0}^{\infty} x^n = \frac{x^2}{2} \frac{1}{1-x} = \frac{3x^2}{4};$$

and hence

$$\frac{1}{\log(1-x)} = \frac{1}{x+3x^2/4} = \frac{1}{x+x/4} = \frac{4}{5x}; \quad (\text{A.12})$$

From the definition of  $T_j$  in (2.2), for  $j \geq 4$ ,  $3j \geq 3$ ,  $1 \geq 3$  we have

$$T_j = \frac{\log \frac{j-4-3j}{B^8}}{\log \frac{1-3j-4}{3j}} = \frac{12}{5j-4-3j} = \frac{1}{2} \log \frac{j-4-3j}{B^8} = \frac{1}{2} \log \frac{j-4-3j}{B^8}; \quad (\text{A.13})$$

Under scaling condition (2.6), along with relation  $\frac{B^4}{j-4-3j} = \frac{1}{8}$  given by Lemma 12 in Appendix B and  $T_j \geq 100$ , we have

$$\frac{(T_j)_{j-4-3j}}{3} = \frac{3 \log \frac{2j-4-3j^2}{9B^8} T_j}{2T_j} < \frac{3 \log T_j}{T_j} = \frac{1}{3}; \quad \frac{2j-4-3j^2}{9B^8} T_j \leq e; \quad (\text{A.14})$$

Plugging in  $T_j = T$  (as in (2.5)) to (A.13) and we obtain

$$\begin{aligned} T_{(T)} &= \frac{12}{5j-4-3j} = \frac{1}{2} \log \frac{j-4-3j}{B^8} = \frac{1}{2} \log \frac{j-4-3j}{B^8} \\ &= \frac{12}{5} \frac{2T}{9 \log \frac{2j-4-3j^2}{9B^8} T} = \log \frac{2j-4-3j^2 T}{9B^8 \log \frac{2j-4-3j^2}{9B^8} T} \\ &= \frac{8}{15} \frac{T}{\log \frac{2j-4-3j^2}{9B^8} T} = \frac{1}{2} \log \frac{2j-4-3j^2}{9B^8} T = \frac{4}{15} T; \end{aligned} \quad (\text{A.15})$$

where we use the elementary inequality  $\log \frac{x}{\log x} \geq \frac{1}{2} \log x$  for all  $x \geq e$ .

(3) From (A.11) and (A.15) we know that  $2 \leq [T_{(T);0:5}; T_{(T);4}]$ . Here we will verify scaling condition (2.3) required in Lemma 3 under our setting. By choosing

$$36 + \frac{5184}{\log^5 1} \leq d; \quad (\text{A.16})$$

we have

$$T_{(T);1} \geq C_{4,T}^0 \leq dT; \quad (\text{A.17})$$

where we define constant  $C_{4,T}^0 = (4-3)(36+5184) = 6960$  and use result  $T_{(T);1} \geq 4T = 3e^{-1}$  implied by (A.11),  $T \geq 1$ .

Therefore for the first scaling condition in (2.3), our pick of  $d = 4$  requires

$$225C_{3,L} \log^8(T_{(T);1}) \leq \frac{d \log \frac{2j-4-3j^2}{9B^8} T}{T} = \log \frac{2j-4-3j^2 T}{9B^8 \log \frac{2j-4-3j^2}{9B^8} T} = \frac{j-4-3j^2}{B^8};$$

while a sufficient condition for the above to hold is, due to (A.17),

$$C_{4;T} \log^8(C_{4;T}^{-1} dT) \frac{d \log^2 T}{T} \leq \frac{j^4 - 3j^2}{B^8}; \quad (\text{A.18})$$

which comes from (2.6) and constant definition of  $C_{4;T} = 225C_{3;L}^2 = 90C_{3;L}$ , because  $T \geq 100$  and the relation (B.1) of  $B; j^4 - 3j^2$  in Appendix Lemma 12 imply  $\frac{j^4 - 3j^2}{9B^8} < \frac{1}{9}$ , and hence

$$1 - \log \frac{j^4 - 3j^2}{9B^8} T \leq 2 \log T; \quad (\text{A.19})$$

To verify the second condition in (2.3), using scaling condition (2.6) in Theorem 4 and (A.19), we have

$$\frac{B^8}{j^4 - 3j^2} (T) = \frac{9B^8}{2j^4 - 3j^2} \frac{\log \frac{j^4 - 3j^2}{9B^8} T}{T} \leq \frac{9B^8}{j^4 - 3j^2} \frac{\log T}{T} < e^{-1};$$

Note that we have already verified  $\frac{j^4 - 3j^2}{9B^8} (T) < 1$  in (A.14). So far we have shown that, under scaling condition (2.6) in Theorem 4, for stepsize  $\epsilon(T)$  given in (2.5) the scaling condition (2.3) in Lemma 3 is guaranteed to hold, which justifies our following act on proving Theorem 4 with Lemma 3.

- (4) Using warm initialization condition (2.1) in Lemma 3 and the definition of  $\epsilon(T)$  given in (2.2), for all  $t \in [T_{0.5}; T_4]$  we have

$$\begin{aligned} \tan \setminus u^{(0)}; a_i &\leq \frac{1}{3} \frac{j^4 - 3j^2}{B^8} \frac{1}{s} \frac{B^4}{j^4 - 3j^2} \frac{1}{3} \\ &\leq \frac{1}{3} \log^{5=2} \frac{1}{s} \frac{B^4}{j^4 - 3j^2} \frac{1}{d \log \frac{j^4 - 3j^2}{B^8}}; \end{aligned}$$

where the first inequality comes from  $T_{0.5}$  and definition of  $T_4$  in (2.2), and the second inequality is due to  $\frac{1}{3} \leq e$  and  $\frac{j^4 - 3j^2}{B^8} \leq 1$  given by scaling condition (2.3). Therefore, on the event  $\mathcal{H}_{3;L}$  we have for all  $t \in [T_{0.5}; T_4]$

$$\tan \setminus u^{(t)}; a_i \leq 3C_{3;L} \log^{5=2} \frac{1}{s} \frac{B^4}{j^4 - 3j^2} \frac{1}{d \log \frac{j^4 - 3j^2}{B^8}}; \quad (\text{A.20})$$

To finalize our proof, we plug in  $\epsilon = \epsilon(T)$  to (A.20) and conclude from  $t \in [T_{(T);0.5}; T_{(T);4}]$  that there exists an event  $\mathcal{H}_{4;T}$  equivalent to  $\mathcal{H}_{3;L}$  with, due to (A.16),  $P(\mathcal{H}_{4;T}) \geq 1 - \frac{1}{3}$ ; such



that on event  $\mathcal{H}_{4,T}$  the following holds

$$\begin{aligned} & \tan \setminus u^{(T)}; a_i \\ & 3C_{3,L} \log^{5=2} \frac{B^4}{j^4 3^{j^2}} \frac{d(T) \log \frac{j^4 3^j}{B^8} (T)^{-1}}{s} \\ & = \frac{9^p}{2} C_{3,L} \log^{5=2} \frac{B^4}{j^4 3^j} \frac{d \log \frac{2j^4 3^{j^2} T}{9B^8} \log \frac{2j^4 3^{j^2} T}{9B^8 \log \frac{2j^4 3^{j^2} T}{9B^8}}}{T} \frac{0}{1} \\ & C_{4,T} \log^{5=2} (C_{4,T}^0 d^{-1}) \frac{B^4}{j^4 3^j} \frac{d \log^2 T}{T}; \end{aligned} \tag{A.21}$$

where in the last step we apply (A.19)  $\log^{5=2} (C_{4,T}^0 d^{-1})$  from (A.16), with constants  $C_{4,T}^0, C_{3,L}, C_{4,T}$ . This completes the whole proof of the theorem.  $\blacksquare$

### A.3. Proof of Lemma 6

In uniform initialization analysis, intuitively  $v^{(t)}$  needs to enter intermediate region  $\mathcal{D}_{\text{mid};k}$  first before we worry about its exit of the intermediate-auxiliary region  $\mathcal{D}_{\text{mid-aux};k}$ . For each coordinate  $k \in [2; d]$ , we upper bound the stopping time  $T_{c;k}$  earlier defined in (2.12) using cold initialization Lemmas 10 and 11. For each coordinate  $k \in [2; d]$ , we view the iterative process  $v_k^{(t)} = v_1^{(t)} g$  as a Markov chain with  $v$  initialized in region  $\mathcal{D}_{\text{mid-aux};k}$  shifted by  $T_{c;k}$ . We notice that process  $U_k^{(t)} 1_{(t < T_1)}$  is bounded by 1 due to definition of  $T_1$ , and we have  $U_k^{(t)} 1_{(t < T_1)} = U_k^{(t)}$  for all  $(t < T_{w;k})$ . Due to strong Markov property the intermediate initialization Lemma 9 applies to the shifted Markov chain.

Proof [Proof of Lemma 6] We let constant  $C_{6,L} = \max\{256C_{11,L}^2; 476C_{9,L}^2\}$  in scaling condition (3.2) in Lemma 6.

- (1) We start by making coordinate-wise analysis for each  $k \in [2; d]$ . For all  $T_{c;k}^0 \wedge T_{c;k} \wedge T_1$ , on the event  $\mathcal{H}_{10,L} \setminus \mathcal{H}_{k;11,L}$ , since  $(v_1^{(t-1)})^2 = 1/d$ , by applying Lemmas 10 and 11 we have

$$\begin{aligned} W_k^{(t)} - W_k^{(0)} & \leq C_{11,L} \log^{5=2} \frac{B^4}{j^4 3^{j^2}} d^{1=2} \left(1 + \frac{1}{d} \frac{j^4 3^j}{j^4 3^j}\right)^t \\ & \leq \frac{C_{11,L} \log^{5=2} \frac{B^4}{j^4 3^{j^2}} d^{1=2} \left(1 + \frac{1}{d} \frac{j^4 3^j}{j^4 3^j}\right)^t}{8 \log^{-1} \log d} \leq 0; \end{aligned} \tag{A.22}$$

where in the last step we used the following inequality implied by scaling condition (3.2)

$$\frac{C_{11,L} \log^{5=2} \frac{B^4}{j^4 3^{j^2}} d^{1=2}}{8 \log^{-1} \log d} \leq 2C_{11,L} \log^{5=2} \frac{B^4}{j^4 3^{j^2}} d^{1=2}; \tag{A.23}$$

Using the elementary inequality  $\log(1-x) \leq \log(1+2x)$  for all  $0 < x < \frac{1}{2}$ , we have

$$1 + \frac{j}{d} \frac{j-4}{4} \frac{3j}{3j} T_{0.5}^{\circ} \exp\left\{\frac{1}{2} \log \frac{j-4}{B^8} \frac{3j}{3j} - \frac{\log\left(1 + \frac{j-4}{d} \frac{3j}{3j}\right)}{\log\left(1 - \frac{j-4}{2d} \frac{3j}{3j}\right)} A \frac{j-4}{B^4} \frac{3j^{1=2}}{B^4}\right\}, \quad (\text{A.24})$$

since  $\frac{j-4}{2d} \frac{3j}{3j} < \frac{1}{2}$  under scaling condition (3.2). Hence on the event

$$H_{10;L} \setminus H_{k;11;L} \setminus (T_1 > T_{0.5}^{\circ}) \setminus (T_{c;k} > T_{0.5}^{\circ});$$

using (A.22), (A.23), (A.24) and  $2 \in (0; e^{-1}]$ , we have

$$W_k^{(T_{0.5}^{\circ})} \leq C_{11;L} \log^{5=2} \frac{1}{j-4} \frac{B^4}{3j^{1=2}} d^{1=2} \frac{j-4}{B^4} \frac{3j^{1=2}}{3j^{1=2}} = C_{11;L} \log^{5=2} \frac{1}{j-4} d^{-2}; \quad (\text{A.25})$$

which indicates  $v^{(T_{0.5}^{\circ})} \in D_{\text{mid-aux}k}$ , i.e. event  $H_{10;L} \setminus H_{k;11;L} \setminus (T_1 > T_{0.5}^{\circ}) \setminus (T_{c;k} > T_{0.5}^{\circ})$ , and hence

$$H_{10;L} \setminus H_{k;11;L} \setminus (T_1 > T_{0.5}^{\circ}) \setminus (T_{c;k} > T_{0.5}^{\circ}); \quad (\text{A.26})$$

- (2) To study  $v^{(t)}$  in the intermediate region, we first assume  $\bar{\tau}_{\theta,k} = 0$ , that is, the initialization  $v^{(0)} \in D_{\text{mid-aux}k}$ . For all  $T_{w;k} \wedge T_{0.5}^{\circ}$ , on the event  $H_{k;9;L}$ , since  $(v_1^{(t-1)})^2 \leq (v_k^{(t-1)})^2 < \frac{1}{2}$  and the following is guaranteed by scaling condition (3.2)

$$\frac{j-4}{2d} \frac{3j}{3j} < \frac{1}{2}; \quad 2C_{9;L} B^4 \log^{5=2} \frac{1}{j-4} d^{1=2} (T_{0.5}^{\circ})^{1=2} \leq \frac{r}{476}; \quad (\text{A.27})$$

applying Lemma 9 we have

$$\begin{aligned} |U_k^{(t)} - U_k^{(0)}| &\leq \frac{j-4}{2d} \frac{3j}{3j} \frac{1}{2} + 2C_{9;L} B^4 \log^{5=2} \frac{1}{j-4} d^{1=2} (T_{0.5}^{\circ})^{1=2} \\ &< \frac{1}{3} \frac{1}{2} + \frac{r}{476} \frac{1}{2}. \end{aligned} \quad (\text{A.28})$$

Now we consider uniform initialization case. Using strong Markov property as discussed earlier, we have for all  $t \geq [T_{c;k}; T_{w;k} \wedge T_{0.5}^{\circ}]$

$$U_k^{(t)} \leq \frac{1}{2}; \quad (\text{A.29})$$

On the event

$$H_{k;9;L} \setminus H_{10;L} \setminus H_{k;11;L} \setminus (T_1 > T_{0.5}^{\circ}) \setminus (T_{w;k} > T_{0.5}^{\circ});$$

we have already proved  $\bar{\tau}_{c;k} > T_{0.5}^{\circ}$  ( $> 0.5$ ) in (A.26). Applying (A.29) with  $t = T_{w;k} \wedge T_{0.5}^{\circ}$ , we obtain  $|U_k^{(T_{w;k})} - U_k^{(0)}| \leq \frac{1}{2}$ , which leads to contradiction with the definition of  $T_{w;k}$  in (3.9) and indicates  $H_{k;9;L} \setminus H_{10;L} \setminus H_{k;11;L} \setminus (T_1 > T_{0.5}^{\circ}) \setminus (T_{w;k} > T_{0.5}^{\circ}) = \emptyset$ , i.e.

$$H_{k;9;L} \setminus H_{10;L} \setminus H_{k;11;L} \setminus (T_1 > T_{0.5}^{\circ}) \setminus (T_{w;k} > T_{0.5}^{\circ}); \quad (\text{A.30})$$

(3) With (A.26) and (A.30) proven, we put all coordinates  $k \in [2; d]$  together and define event

$$H_{6;L} = \bigcap_{k \in [2; d]} H_{k;9;L} \cap H_{10;L} \cap \bigcap_{k \in [2; d]} H_{k;11;L} :$$

On the event  $H_{6;L} \setminus (T_1 \leq T_{0;0.5}^0)$ , for each coordinate  $k \in [2; d]$  satisfying  $T_1 \leq T_{0;0.5}^0 \wedge T_{c;k}$ , we apply (A.22) with  $t = T_1 = T_{0;0.5}^0 \wedge T_{c;k} \wedge T_1$  and obtain  $W_k^{(T_1)} = 0$ . Then due to the definition of  $T_1$  in (3.15), there must exist  $k \in [2; d]$  such that  $W_k^{(T_1)} < 0$  and  $T_{c;k} < T_1 \leq T_{0;0.5}^0$ . Recall that  $W_k^{(t)} < 0$  is equivalent to  $j U_k^{(t)} > 1$ . By definitions of stopping times  $T_1, T_{w;k}$ , we know the existence of  $k \in [2; d]$  such that  $T_{c;k} < T_{w;k} \leq T_1 \leq T_{0;0.5}^0$  on the event  $H_{6;L} \setminus (T_1 \leq T_{0;0.5}^0)$ . By applying (A.29) with  $t = T_{w;k} = T_{w;k} \wedge T_{0;0.5}^0$ , we find  $j U_k^{(T_{w;k})} \leq \frac{1}{2}$ , which contradicts with the definition of  $T_{w;k}$  and indicates that  $H_{6;L} \setminus (T_1 \leq T_{0;0.5}^0) = \emptyset$ , i.e.

$$H_{6;L} \subseteq (T_1 > T_{0;0.5}^0): \quad (\text{A.31})$$

Combining (A.26) and (A.30) for each  $k \in [2; d]$  along with (A.31), we have

$$H_{6;L} \subseteq \left( \sup_{k \in [2; d]} T_{c;k} \leq T_{0;0.5}^0 \right) \cap \left( \inf_{k \in [2; d]} T_{w;k} > T_{0;0.5}^0 \right) \cap (T_1 > T_{0;0.5}^0): \quad (\text{A.32})$$

(4) With (A.32) ready at hand, on the event  $H_{6;L}$ , we apply Lemma 9 for each coordinate  $k \in [2; d]$  and all  $t \in [T_{0;0.5}^0, T_{0;0.5}^0]$  to obtain

$$\begin{aligned} & \sum_{k=2}^d \sum_{t=T_{0;0.5}^0}^0 \left\| \sum_{s=T_{0;0.5}^0}^1 \sum_{j=4}^3 \sum_{i=1}^2 \frac{X_k^{(t)} \otimes U_k^{(T_{0;0.5}^0)} \otimes \mathbf{1}}{(v_1^{(s)})^2 (v_k^{(s)})^2} \right\|_A^2 \\ & \leq 2C_{9;L} \log^{5-2} \frac{1}{B^4} d (T_{0;0.5}^0)^{1-2}. \end{aligned} \quad (\text{A.33})$$

On the event  $H_{6;L}$ , we have  $T_{c;k} \leq T_{0;0.5}^0, (v_1^{(s)})^2 (v_k^{(s)})^2 \leq 1/(2d), U_k^{(T_{c;k})} \leq 1$  for all  $k \in [2; d], s \in [T_{0;0.5}^0, T_{0;0.5}^0]$ . Since the left hand of (A.33) is the norm of two vectors subtraction, we use the triangle inequality of Euclidean norms to lower bound it as

$$\begin{aligned} & \sum_{k=2}^d \sum_{t=T_{0;0.5}^0}^0 \left\| \sum_{s=T_{0;0.5}^0}^1 \sum_{j=4}^3 \sum_{i=1}^2 \frac{X_k^{(t)} \otimes U_k^{(T_{0;0.5}^0)} \otimes \mathbf{1}}{(v_1^{(s)})^2 (v_k^{(s)})^2} \right\|_A^2 \\ & \geq \sum_{k=2}^d \sum_{t=T_{0;0.5}^0}^0 \left\| \sum_{s=T_{0;0.5}^0}^1 \sum_{j=4}^3 \sum_{i=1}^2 \frac{X_k^{(t)} \otimes U_k^{(T_{0;0.5}^0)} \otimes \mathbf{1}}{(v_1^{(s)})^2 (v_k^{(s)})^2} \right\|_A^2 \\ & \geq \tan \left( \frac{1}{2d} \right) \sum_{k=2}^d \sum_{t=T_{0;0.5}^0}^0 \left\| \sum_{s=T_{0;0.5}^0}^1 \sum_{j=4}^3 \sum_{i=1}^2 \frac{X_k^{(t)} \otimes U_k^{(T_{0;0.5}^0)} \otimes \mathbf{1}}{(v_1^{(s)})^2 (v_k^{(s)})^2} \right\|_A^2. \end{aligned} \quad (\text{A.34})$$

Recall the definition of  $T_{0;0.5}^0$  in (3.1). Scaling condition (3.2) guarantees the following

$$\frac{j=4}{2d} \sum_{j=4}^3 \frac{1}{3j} \geq 1; \quad \frac{B^8}{j=4 \sum_{j=4}^3} \geq e^{-1}: \quad (\text{A.35})$$

Since  $\log \frac{1}{2d} \frac{j^4 - 3j}{j^4 - 3j}$  and (A.35) holds, for each positive we have relation

$$T_{;}^{\circ} = 1 + \frac{2d}{j^4 - 3j} \log \frac{j^4 - 3j}{B^8} = \frac{2(+1)d}{j^4 - 3j} \log \frac{j^4 - 3j}{B^8} \quad (A.36)$$

Combining (A.33), (A.34) together and using relation (A.36), we have

$$\begin{aligned} \tan \setminus v^{(t)}; e_1 &= \frac{p}{d} \frac{1}{2d} \frac{j^4 - 3j}{3j} t T_{;0;5}^{\circ} \\ &+ \frac{p}{+1} C_{6;L} \log^{5=2} \frac{B^4}{j^4 - 3j^{1=2}} d^3 \log \frac{j^4 - 3j}{B^8} \quad ; \end{aligned} \quad (A.37)$$

where constant  $C_{6;L} = 2^p C_{9;L}$ . To complete proof of Lemma 6, we provide a lower bound on probability of event  $H_{6;L}$  by taking union bound

$$P(H_{6;L}) = \prod_{k=2}^X P(H_{k;9;L}^c) P(H_{10;L}^c) \prod_{k=2}^X P(H_{k;11;L}^c) = 1 - 6 + 27 + \frac{10368}{\log^5 1} d^{-3} :$$

Applying the scaling relation (2.9) of  $u^{(t)}_{g_t=0}$  and  $v^{(t)}_{g_t=0}$  to (A.37) on the event  $H$  completes the proof of (3.3), and hence Lemma 6. ■

#### A.4. Proof of Theorem 7

Now we are ready to derive the finite-sample error bound and prove Theorem 7.

Proof [Proof of Theorem 7]

- (1) Analogous to the proof of Theorem 4, we sharply upper- and lower-bound the rescaled time  $T_{(T)}^{\circ}$ . Using relation  $\frac{B^4}{j^4 - 3j} = \frac{1}{8}$  in Lemma 12 and elementary inequality  $\frac{\log T}{T} \leq \frac{1}{20}$ , we have

$$\frac{j^4 - 3j}{2d} (T) = \frac{2}{T} \log \frac{j^4 - 3j^2}{4B^8 d} T = \frac{1}{3} \quad (A.38)$$

Since  $\frac{1}{\log(1-x)} \leq \frac{1}{x}$  for  $x \in (0, 1)$ , similar to (A.36) we have

$$\begin{aligned} T_{(T)}^{\circ} &= 1 + \frac{2d}{j^4 - 3j} (T) \log \frac{j^4 - 3j}{B^8} (T) \\ &= 1 + \frac{T \log \frac{j^4 - 3j^2}{4B^8 d} T}{2 \log \frac{j^4 - 3j^2}{4B^8 d} T} \log \frac{j^4 - 3j^2}{4B^8 d} T \\ &= 1 + \frac{T}{2} \frac{\log \frac{j^4 - 3j^2}{4B^8 d} T}{\log \frac{j^4 - 3j^2}{4B^8 d} T} + \frac{1}{2} T \quad ; \end{aligned} \quad (A.39)$$

where the last inequality holds due to  $\frac{j^4 - 3j^2}{4B^8 d} T$  under scaling condition (3.2) with constants  $C_{7;T} = 96C_{6;L}$ ;  $C_{7;T} = 10425$ .

On the lower bound side, previously in (A.12) we show that  $\frac{1}{\log(1-x)} \leq \frac{4}{5x}$  for all  $x \in (0; 1/3]$ , which can be applied to  $T_{(T);1}^0$  by replacing  $x$  with  $\frac{j-4}{2d} \frac{3j}{T}$  (T) due to (A.38). By noticing that  $T \geq 100$  and scaling condition (3.4) guarantee the following inequalities

$$\frac{4 \log T}{T} \leq \frac{1}{3}; \quad \frac{8B^8}{j-4} \frac{d \log T}{3j^2 T} \leq 1; \quad (\text{A.40})$$

we have

$$T_{(T);1}^0 = \frac{\log \frac{j-4}{B^8} \frac{3j}{T}}{\log \left(1 - \frac{j-4}{2d} \frac{3j}{T}\right)} \leq \frac{8 d \log \frac{j-4}{B^8} \frac{3j}{T}}{5j-4} \frac{3j}{T} \leq \frac{2T}{5 \log \frac{j-4}{4B^8 d} T} \log \left( \frac{j-4}{4B^8 d} \frac{3j^2 T}{4B^8 d} \right) \leq \frac{1}{5} T; \quad (\text{A.41})$$

where in the last step we use the elementary inequality  $\frac{x}{\log x} \leq \frac{1}{2} \log x$  for all  $x \geq e$ , since  $\frac{j-4}{4B^8 d} T \geq e$  is satisfied under scaling condition (3.2).

- (2) From (A.39) and (A.41) we find that  $T \leq 2 [T_{(T);1}^0 + 1; T_{(T);5}^0]$ . By letting  $\alpha = 57 + \frac{10368}{\log^5 1} d$  along with (A.39) we have

$$T_{(T);1}^0 \leq C_{7;T}^0 \frac{1}{d} T;$$

for positive, absolute constant  $C_{7;T}^0 \leq 10425$  due to  $T_{(T);1}^0 \leq T$  given by (A.39) and  $e^{-1}$  guaranteed by  $\alpha \geq 4$ .

The third scaling condition in (3.2) with our pick  $\alpha = 5$  and  $\beta = (T)$  is satisfied by

$$24C_{6;L} \log^8(C_{7;T}^0 \frac{1}{d} T) \leq \frac{B^8}{j-4} \frac{3j^2}{T} \frac{d^3 \log^2 d \log \frac{j-4}{4B^8 d} T}{T} \log \left( \frac{j-4}{4B^8 d} \frac{3j^2 T}{4B^8 d} \right) \leq \frac{1}{\log^2 1}; \quad (\text{A.42})$$

From Lemma 12 we have  $\frac{B^4}{j-4} \frac{3j}{T} \leq \frac{1}{8}$ . Along with scaling condition (3.4) and  $C_{7;T}^0 \leq 96C_{6;L}$ , we have  $T \geq d \geq 16$  and hence

$$1 \leq \log \frac{j-4}{4B^8 d} \frac{3j^2 T}{4B^8 d} \leq 2 \log(T=d); \quad (\text{A.43})$$

implying (A.42) holds under scaling condition (3.4).

To verify the second scaling condition in (3.2), we notice that the following holds under (3.4)

$$\frac{B^8}{j-4} \frac{3j}{T} = \frac{4B^8 d \log \frac{j-4}{4B^8 d} T}{j-4} \frac{3j^2 T}{3j^2 T} \leq \frac{8B^8 d \log T}{j-4} \frac{3j^2 T}{3j^2 T} < e^{-1};$$

and

$$(T) = \frac{4d \log \frac{j-4}{4B^8 d} T}{j-4} \frac{3j^2 T}{3j^2 T} \leq \frac{8 \log(T=d)}{j-4} \frac{3j^2 T}{3j^2 T} < \frac{1}{j-4} \frac{3j^2 T}{3j^2 T};$$

due to (A.43) and the elementary inequality  $\frac{8 \log x}{x} < 1$  for all  $x \geq 100$ .

Therefore, all scaling conditions required in Lemma 6 are satisfied by scaling condition (3.4) in Theorem 7.

- (3) Using  $B^8 = j^4 3^j e^{-1}$  and  $e^{-1}$  given by (3.2) and  $T_{;1}^0 + 1 \leq T_{;0:5}^0 \leq T_{;0:5}^0$  following its definition (3.1), for all  $2 \in [T_{;1}^0 + 1; T_{;5}^0]$ , on the event  $H_{6;L}$  we have

$$\begin{aligned} & \frac{p_{\bar{d}} - 1}{2^d} j^4 3^j e^{-T_{;0:5}^0} \leq \frac{p_{\bar{d}} - 1}{2^d} j^4 3^j e^{-T_{;0:5}^0} \frac{B^4}{j^4 3^{j^{1=2}}} \frac{p_{\bar{d}}}{d} \\ & (3.6) C_{6;L} \log^{5=2} \frac{B^4}{j^4 3^{j^{1=2}}} \leq \frac{B^4}{d^3 \log \frac{j^4 3^j}{B^8}} \frac{1}{d} : \end{aligned}$$

Therefore, from Lemma 6, on the event  $H_{6;L}$  we have for all  $2 \in [T_{;1}^0 + 1; T_{;5}^0]$  that

$$\tan \setminus u^{(t)}; a_l \leq 3 C_{6;L} \log^{5=2} \frac{B^4}{j^4 3^{j^{1=2}}} \frac{1}{d^3 \log \frac{j^4 3^j}{B^8}} \frac{1}{d} : \quad (\text{A.44})$$

Plugging in our choice of  $d = 5$  and  $(T) = \frac{4d \log \frac{j^4 3^j T}{4B^8 d}}{j^4 3^j T}$  to (A.44) and using (A.43), we know that there exists an event  $H_{7;T} \subset H_{6;L}$  with  $P(H_{7;T}) \geq 1 - 4^{-4}$  such that on event  $H_{7;T}$  we have

$$\tan \setminus u^{(T)}; a_l \leq C_{7;T} \log^{5=2} (C_{7;T}^0 \frac{1}{d}) \frac{B^4}{j^4 3^j} \frac{1}{\frac{d^4 \log^2 T}{T}};$$

where constant  $C_{7;T} = 12 C_{6;L}; C_{7;T}^0 = 10425$  ■

## Appendix B. Secondary Lemmas in Warm Initialization Analysis

For notational simplicity, we denote  $v^{(t-1)}$  and  $Y = Y^{(t)}$ . We first provide a lemma on Orlicz  $2$ -norm of  $v^{(t-1)}$  and the relation between  $B$  and  $d$ .

**Lemma 12** Let Assumption 2 hold. For each rotated observation  $v_k$  and any unit vector  $w$ , we have Orlicz  $2$ -norm  $\|v_k\|_2 \leq B$  and the following relation of  $B$  and  $d$

$$\frac{B^4}{j^4 3^j} \leq \frac{1}{8} \quad (\text{B.1})$$

With the bound on Orlicz  $2$ -norm given in Lemma 12 and  $\bar{d}_{;1}$  defined in (2.2), we introduce truncation barrier parameter

$$B = B \log^{1=2}(T_{;1} - 1); \quad (\text{B.2})$$

where  $\epsilon \in (0; e^{-1}]$  is some fixed positive. For each coordinate  $k \in [2; d]$  define the first time the norm of a data observation exceeds the truncation barrier as

$$T_{B;k} = \inf_{t \in [1; n]} \{ v^{(t-1)} > Y^{(t)} > B \text{ or } Y_1^{(t)} > B \text{ or } Y_k^{(t)} > B \}; \quad (\text{B.3})$$

B.1. Proof of Lemma 5

For each  $k \in [1]$ , we define random variable

$$Q_{U;k}^{(t)} = U_k^{(t)} - U_k^{(t-1)} - \text{sign}(v_k - Y_k) (v_k - Y_k)^3 v_1^2 (v_1 Y_k - v_k Y_1); \quad (\text{B.4})$$

Lemma 13 Let  $B$  be any positive value. For each coordinate  $k \in [2; d]$  and any  $t \geq 1$ , on the event

$$H_{k;13;L}^{(t)} = \{ |v_k - Y_k| \leq B; |Y_1| \leq B; |Y_k| \leq B; v_1^2 \leq d; v_1^2 \leq v_k^2; \} \quad (\text{B.5})$$

for stepsize  $\eta = (2B^4 d^{1/2})^{-1}$  we have  $\mathbb{E} |Q_{U;k}^{(t)}| \leq 4B^8 \eta^2 v_1^2$ .

Lemma 14 Let Assumption 2 hold. For each coordinate  $k \in [2; d]$  and any  $t \geq 1$ , we have

$$\mathbb{E} \sum_{k=1}^h \text{sign}(v_k - Y_k) (v_k - Y_k)^3 v_1^2 (v_1 Y_k - v_k Y_1) F_{t-1} = \sum_{k=1}^h \sum_{j=1}^3 (v_1^2 - v_k^2) U_k^{(t-1)}; \quad (\text{B.6})$$

For each  $k \in [2; d]$  and  $t \geq 1$ , at the  $t$ -th iteration we define

$$e_k^{(t)} = \text{sign}(v_k - Y_k) (v_k - Y_k)^3 v_1^2 (v_1 Y_k - v_k Y_1) + \sum_{j=1}^3 (v_1^2 - v_k^2) U_k^{(t-1)} \quad (\text{B.7})$$

which, indexed by  $t$ , forms a sequence of martingale differences with respect to  $\mathcal{F}_t$ .

Lemma 15 For each coordinate  $k \in [2; d]$  and any  $t \geq 1$ ,  $U_k^{(t)}$  has linear representation

$$U_k^{(t)} = U_k^{(0)} + \sum_{s=0}^{t-1} \sum_{j=1}^3 (v_1^{(s)})^2 - (v_k^{(s)})^2 U_k^{(s)} + \sum_{s=1}^t Q_{U;k}^{(s)} + \sum_{s=1}^t e_k^{(s)}; \quad (\text{B.8})$$

Lemma 16 Let Assumption 2 hold and initialization  $U^{(0)} \in D_{\text{warm}}$ . Let  $\epsilon \in (0; e^{-1}]$  and  $\delta$  be any fixed positive. For each coordinate  $k \in [2; d]$ , there exists an event  $H_{k;16;L}$  satisfying

$$P(H_{k;16;L}) \geq 1 - \delta + \frac{5184}{\log^5 \frac{1}{\delta}};$$

such that on the event  $H_{k;16;L}$  the following concentration result holds

$$\max_{1 \leq t \leq T; \forall x} \sum_{s=1}^t e_k^{(s)} \leq C_{16;L} \log^{\frac{5}{2}} \frac{1}{\delta} B^4 (T; \epsilon)^{1/2};$$

where  $C_{16;L}$  is a positive, absolute constant.

Lemma 17 Let  $\epsilon \in (0; e^{-1}]$  and  $\delta$  be any fixed positive. For each coordinate  $k \in [2; d]$ , we have

$$P(T_{B;\delta;k} \leq T; \epsilon) \geq 6(1 + \delta); \quad (\text{B.9})$$

With the above secondary lemmas at hand, we are now ready to prove Lemma 5. Proof [Proof of Lemma 5] We recall the definition of stopping time  $T_{k;5;L}$  in (2.12). Because stepsize  $1=(2B^4d^{1=2})$  holds under scaling condition (2.3), on the event  $(T_{B;k} > T_{; } \wedge T_x) \setminus (T_{B;k} > T_{; } ) \cap H_{k;13;L}^{(t)}$ , we have  $v_1^2 \leq \frac{3}{2}$ , and applying Lemma 13 gives

$$|Q_{U;k}^{(t)}| \leq 4B^8 v_1^2 + 6B^8 \log^4(T_{; } - 1):$$

We define event  $H_{k;5;L} = (T_{B;k} > T_{; } ) \setminus H_{k;16;L}$ . Applying Lemmas 15 and 16, on the event  $H_{k;5;L}$  for all  $t \leq T_{; } \wedge T_x$  we have

$$U_k^{(t)} \leq U_k^{(0)} + \sum_{j=4}^{\lfloor t \rfloor} 3 \sum_{s=0}^{\lfloor t-j \rfloor} (v_1^{(s)})^2 + (v_k^{(s)})^2 U_k^{(s)} \tag{B.10}$$

$$T_{; } \leq 6B^8 \log^4(T_{; } - 1) + C_{16;L} \log^{5=2} T_{; } B^4 (T_{; } - 1)^{1=2}:$$

Scaling condition (2.3) and definition of  $T_{; }$  in (2.2) imply that

$$B^4 (T_{; } - 1)^{1=2} \log^4(T_{; } - 1) \leq \log^{5=2} T_{; } \tag{B.11}$$

Combining (B.10) and (B.11), we have

$$U_k^{(t)} \leq U_k^{(0)} + \sum_{j=4}^{\lfloor t \rfloor} 3 \sum_{s=0}^{\lfloor t-j \rfloor} (v_1^{(s)})^2 + (v_k^{(s)})^2 U_k^{(s)} + C_{5;L} \log^{5=2} T_{; } B^4 (T_{; } - 1)^{1=2};$$

where constant  $C_{5;L} = C_{16;L} + 6$ . Scaling condition (2.3) implies  $\sum_{j=4}^{\lfloor t \rfloor} 3 \sum_{s=0}^{\lfloor t-j \rfloor} ((v_1^{(s)})^2 + (v_k^{(s)})^2) \leq 2$  for all  $s < T_{; } \wedge T_x$ . Using reversed Gronwall Lemma 25, on the event  $H_{k;5;L}$  we have the following holds for all  $t \leq T_{; } \wedge T_x$

$$U_k^{(t)} \leq U_k^{(0)} e^{\sum_{j=4}^{\lfloor t \rfloor} 3 \sum_{s=0}^{\lfloor t-j \rfloor} (v_1^{(s)})^2 + (v_k^{(s)})^2} + 2C_{5;L} \log^{5=2} T_{; } B^4 (T_{; } - 1)^{1=2}:$$

To complete proof of Lemma 5, we apply Lemmas 16 and 17 and take union bound to obtain

$$P(H_{k;5;L}) \leq P(T_{B;k} > T_{; } ) + P(H_{k;16;L}^c) \leq 6 + 12 + \frac{5184}{\log^5 T_{; } - 1} : \quad \blacksquare$$

## B.2. Proof of Secondary Lemmas

Proof [Proof of Lemma 12]

- (1) Because random vector  $Y$  is a permutation of  $Z$ , from Assumption 2 we know that each  $Y_i$  is sub-Gaussian with parameter  $\frac{3}{8}B$ . By independence of  $Y_i$ , for any unit vector  $v$  and all  $t \in \mathbb{R}$  we have

$$E \exp \langle v, Y \rangle^t = \prod_{i=1}^d E \exp \langle v_i, Y_i \rangle^t = \prod_{i=1}^d \exp \left( -\frac{t^2 v_i^2}{2} \right) \leq \left( \frac{3}{8}B \right)^d \exp \left( -\frac{t^2}{2} \right) \leq \left( \frac{3}{8}B \right)^d e^{-\frac{t^2}{2}};$$



which implies that  $v^> Y$  is also sub-Gaussian with parameter  $\sqrt{3-8B}$ . Theorem 2.6 in [Wainwright \(2019\)](#) shows that any sub-Gaussian random variable with parameter  $\sigma$  satisfies  $E \exp \frac{X^2}{2\sigma^2} \leq \frac{1}{1-\epsilon}$  for all  $\epsilon \in [0, 1)$ . By choosing  $\sigma = \sqrt{3-8B}$  and  $X = v^> Y$ , we have

$$E \exp \frac{(v^> Y)^2}{B^2} \leq \frac{1}{1-\epsilon} \quad \text{i.e. } \|v^> Y\|_2 \leq B \quad (\text{B.12})$$

We refer the readers to [Appendix E](#) for more details on Orlicz<sub>2</sub>-norm.

(2) Applying Markov's inequality to [\(B.12\)](#) with  $v = e_i$  gives

$$P(Y_i^2 \geq t) = P \left( \exp \frac{Y_i^2}{B^2} \geq \exp \frac{t}{B^2} \right) \leq \exp \left( -\frac{t}{B^2} \right) E \exp \frac{Y_i^2}{B^2} \leq 2 \exp \left( -\frac{t}{B^2} \right)$$

Hence we have

$$\mathbb{E} Y_i^4 = \int_0^\infty P(Y_i^2 \geq t) dt \leq \int_0^\infty 2 \exp \left( -\frac{t}{B^2} \right) dt = 4B^4;$$

and hence  $\frac{B^4}{j^4} \leq \frac{4}{4j^4} \leq \frac{1}{8}$  holds due to  $\mathbb{E} Y_i^4 = (\mathbb{E} Y_i^2)^2 = 1$ . ■

**Proof [Proof of Lemma 13]** Let  $s = \text{sign}(v_4 - 3)$  for notational simplicity. We recall the definitions of iteration  $U_k^{(t)}$  in [\(2.11\)](#) and random variable  $Q_{U;k}^{(t)}$  in [\(B.4\)](#). Using update formula [\(2.8\)](#), we have

$$\begin{aligned} U_k^{(t)} - U_k^{(t-1)} &= \frac{v_k + s (v^> Y)^3 Y_k}{v_1 + s (v^> Y)^3 Y_1} - \frac{v_k}{v_1} = \frac{v_k + s (v^> Y)^3 Y_k - v_k (v_1 + s (v^> Y)^3 Y_1) / v_1}{(v_1 + s (v^> Y)^3 Y_1) v_1} \\ &= s \frac{1 + s (v^> Y)^3 \frac{Y_1}{v_1} - 1}{v_1} (v^> Y)^3 v_1^2 (v_1 Y_k - v_k Y_1); \end{aligned}$$

Along with [\(B.4\)](#), we obtain

$$Q_{U;k}^{(t)} = s \frac{1 + s (v^> Y)^3 \frac{Y_1}{v_1} - 1}{v_1} (v^> Y)^3 v_1^2 (v_1 Y_k - v_k Y_1);$$

For any  $|x| \leq \frac{1}{2}$ , summation of geometric series gives

$$(1+x)^{-1} - 1 = - \sum_{k=0}^\infty (-x)^k = \sum_{k=0}^\infty (-1)^{k+1} x^k;$$

On the even  $H_{k;13;L}^{(t)}$  defined earlier in [\(B.5\)](#), since  $s (v^> Y)^3 Y_1 = v_1$ ,  $B^4 d^{1-2} = 1=2$ , we have

$$\begin{aligned} |j Q_{U;k}^{(t)}| &\leq \left| 1 + s (v^> Y)^3 \frac{Y_1}{v_1} - 1 \right| |j v^> Y j^3 v_1^2 j v_1 Y_k - v_k Y_1 j| \\ &\leq 2 \left| s (v^> Y)^3 \frac{Y_1}{v_1} \right| |j v^> Y j^3 v_1^2 j v_1 Y_k - v_k Y_1 j| \\ &= 2^2 v_1^2 |j Y_1 j| |j v^> Y j^6 Y_k| \frac{v_k Y_1}{v_1} \leq 4B^8 2 v_1^2. \end{aligned}$$

■

Proof [Proof of Lemma 14] Under Assumption 2, for all  $d \in \mathbb{N}$  we have

$$E \sum_{i=1}^h (v^{\>Y})^3 Y_k F_{t-1}^i = \sum_{i=1}^h v_k^3 + 3v_k \sum_{j \in \mathbb{K}} v_j^2 = (4-3)v_k^3 + 3v_k: \quad (\text{B.13})$$

Recall that  $U_k^{(t-1)} = v_k = v_1$ , then

$$\begin{aligned} & E \sum_{i=1}^h \text{sign}(4-3) (v^{\>Y})^3 v_1^2 (v_1 Y_k - v_k Y_1) F_{t-1}^i \\ &= \text{sign}(4-3) v_1^2 (4-3)v_1 v_k^3 + 3v_1 v_k (4-3)v_k v_1^3 - 3v_k v_1 \\ &= \sum_{j \in \mathbb{K}} (v_1^2 - v_k^2) U_k^{(t-1)}: \end{aligned}$$

■

Proof [Proof of Lemma 15] From definitions (B.4) and (B.7), by applying (B.6) in Lemma 14,

$$U_k^{(s)} - U_k^{(s-1)} = \sum_{j \in \mathbb{K}} (v_1^{(s-1)})^2 - (v_k^{(s-1)})^2 U_k^{(s-1)} + Q_{U,k}^{(s)} + e_k^{(s)}: \quad (\text{B.14})$$

Iteratively applying (B.14) for  $s = 1, \dots, t$  gives (B.8). ■

Proof [Proof of Lemma 16] Under Assumption 2, we apply (E.1) along with Lemmas 12 and 27 to obtain

$$\begin{aligned} & v_1^2 (v^{\>Y})^3 (v_1 Y_k - v_k Y_1) \sum_{i=2}^h v_1^2 k (v^{\>Y})^2 k_{i-1} k (v^{\>Y}) (v_1 Y_k - v_k Y_1) k_{i-1} \\ & v_1^2 k v^{\>Y} k_{i-2}^3 k v_1 Y_k - v_k Y_1 k_{i-2} \\ & \sum_{j \in \mathbb{K}} v_j^{-1} k v^{\>Y} k_{i-2}^3 (k Y_k k_{i-2} + j v_k = v_j k Y_1 k_{i-2}) B^4 \sum_{j \in \mathbb{K}} v_j^{-1} (1 + j v_k = v_j): \end{aligned}$$

Recall the definition of martingale difference sequence  $e_k^{(t)}$  in (B.7). By applying Lemma 28, we have

$$\begin{aligned} e_k^{(t)} &= \sum_{i=2}^h v_1^2 (v^{\>Y})^3 (v_1 Y_k - v_k Y_1) E \sum_{i=1}^h v_1^2 (v^{\>Y})^3 (v_1 Y_k - v_k Y_1) \\ C_{28;L}^0 & \sum_{i=2}^h v_1^2 (v^{\>Y})^3 (v_1 Y_k - v_k Y_1) C_{28;L}^0 B^4 \sum_{j \in \mathbb{K}} v_j^{-1} (1 + j v_k = v_j): \end{aligned} \quad (\text{B.15})$$

Because initialization  $(0) \in D_{\text{warm}}$ , on the event  $(t \leq T_x)$  for  $T_x$  earlier defined in (2.12), we have  $\sum_{j \in \mathbb{K}} v_j^{-1} \leq \frac{1}{3} \leq \frac{1}{2}$ ,  $\sum_{j \in \mathbb{K}} v_j = v_1 \leq \frac{1}{2}$ , and then  $e_k^{(t)} \leq 3C_{28;L}^0 B^4$ . Since  $1_{(t \leq T_x)} \in \mathcal{F}_{t-1}$ , we know that  $e_k^{(t)} 1_{(t \leq T_x)}$  forms a martingale difference sequence with respect to  $\mathcal{F}_t$ . Additionally, because  $e_k^{(t)} 1_{(t \leq T_x)} \leq k e_k^{(t)} k_{i=2}$ , we have

$$e_k^{(t)} 1_{(t \leq T_x)} \leq 3C_{28;L}^0 B^4:$$

With the bound given above, we apply Theorem 29 with  $\delta = 2$  and obtain for all  $\epsilon \in (0; e^{-1}]$ ,

$$\begin{aligned} & P \max_{1 \leq t \leq T; \hat{\Lambda}_{T,x}} \sum_{s=1}^X e_k^{(s)} C_{16;L} B^4 (T; \epsilon)^{1=2} \log^{5=2} \epsilon^{-1} \\ &= P \max_{1 \leq t \leq T; \hat{\Lambda}_{T,x}} \sum_{s=1}^X e_k^{(s)} 1_{(s \neq T_x)} C_{16;L} B^4 (T; \epsilon)^{1=2} \log^{5=2} \epsilon^{-1} \\ & \leq 2 \cdot 3 + 6^4 \frac{64 T; \epsilon}{C_{16;L}^2 B^8 2T; \epsilon} \frac{9 C_{28;L}^2 B^8 2}{\log^5 \epsilon^{-1}} \exp \left\{ \frac{8}{32 T; \epsilon} \frac{C_{16;L}^2 B^8 2T; \epsilon}{9 C_{28;L}^2 B^8 2} \log^5 \epsilon^{-1} \right\}^{\frac{9}{5}} \\ &= 6 + \frac{5184}{\log^5 \epsilon^{-1}} ; \end{aligned}$$

where constant  $C_{16;L} = 12^P \bar{C}_{28;L}^0$ . ■

Proof [Proof of Lemma 17] Recall definition of  $B$  in (B.2),  $T_{B;k}$  in (B.3) and  $T_{;1}$  in (2.2). Using Markov inequality and Lemma 12, we have

$$P(jv^{>Y} j > B) = P \left( \frac{jv^{>Y} j^2}{B^2} > \frac{B^2}{B^2} \right) \leq \exp \left\{ - \frac{B^2}{B^2} \right\} E \exp \left\{ \frac{jv^{>Y} j^2}{B^2} \right\} \leq \frac{2}{T_{;1}}$$

Similarly we also have

$$P(jY_{1j} > B) \leq \frac{2}{T_{;1}}; \quad P(jY_{kj} > B) \leq \frac{2}{T_{;1}}$$

Taking union bound,

$$\begin{aligned} & P_{T_{B;k}} \sum_{t=1}^X P(jv^{(t-1)} > Y^{(t)} j > B) + P(jY_1^{(t)} j > B) + P(jY_k^{(t)} j > B) \\ & \leq 3T_{;1} \frac{2}{T_{;1}} = 6(\epsilon + 1); \end{aligned}$$

where we use the elementary inequality  $\int_0^x (x-t) dt = \frac{x^2}{2}$  for all  $x \geq 0$  to obtain  $T_{;1} = T_{;1} + 1$ . ■

### Appendix C. Secondary Lemmas in Uniform Initialization Analysis

For notational simplicity, we denote  $v^{(t-1)}, Y^{(t)}$ . Recall that we bounded the Orlicz  $\psi_2$ -norm of each  $Y_i$  and  $v^{>Y}$  by  $B$  using Lemma 12 in Appendix B and introduced truncation barrier  $B$  in (B.2) under warm initialization condition, based on rescaled time. Under uniform initialization condition, we consider a different rescaled time defined in (3.1). Accordingly, we introduce a slightly larger truncation barrier based on

$$B_o = B \log^{1=2}(T_{;1}^o); \tag{C.1}$$

where  $\epsilon \in (0; e^{-1}]$  is some fixed positive. For each coordinate  $k \in [2; d]$  define the first time the norm of a data observation exceeds the truncation barrier

$$T_{B_o;k} = \inf_{t \geq 1} \{ jv^{(t-1)} > Y^{(t)} j > B_o \text{ or } jY_1^{(t)} j > B_o \text{ or } jY_k^{(t)} j > B_o \}; \tag{C.2}$$

## C.1. Proof of Lemma 9

Proof of Lemma 9 shares Lemma 13, 14 and 15 with proof of Lemma 5. With the new truncation barrier  $B_0$  given in (C.1) and the new rescaled time  $\tau_0$  earlier defined in (3.1), we plug  $B = B_0$  in Lemma 13 and introduce a new concentration lemma and a new tail probability lemma.

Lemma 18 For any fixed coordinate  $k \in [2; d]$ , let Assumption 2 hold and initialization  $U^{(0)} \in \mathcal{D}_{\text{mid},k}$ . Let  $\epsilon$  be any fixed positives. Then there exists an event  $\mathcal{H}_{k;18;L}$  satisfying

$$P(\mathcal{H}_{k;18;L}) \geq 1 - 6 + \frac{5184}{\log^{5-2} \epsilon};$$

such that on event  $\mathcal{H}_{k;18;L}$  the following concentration result holds

$$\max_{1 \leq t \leq T_0^* \wedge T_{w;k}} \sum_{s=1}^t e_k^{(s)} \leq C_{18;L} \log^{5-2} \epsilon^{-1} B^4 d^{1-2} (T_0^*)^{1-2};$$

where  $C_{18;L}$  is a positive, absolute constant.

Lemma 19 Let  $\epsilon$  be any fixed positive. For each coordinate  $k \in [2; d]$ , we have

$$P(T_{B_0;k} \leq T_0^*) \geq 6(1 - \epsilon). \quad (\text{C.3})$$

With the above secondary lemmas at hand, we are now ready for the proof of Lemma 9. Proof [Proof of Lemma 9] Recall the definition of stopping time  $T_{w;k}$  in (3.9) and  $T_{B_0;k}$  in (C.2). Since scaling condition (3.2) implies stepsize  $\epsilon = 1/(2B_0^4 d^{1-2})$ , we can apply Lemma 13 with  $B = B_0$ . On the event  $(t \leq T_0^* \wedge T_{w;k}) \setminus (T_{B_0;k} > T_0^*)$ , we have  $v_1^2 \leq \epsilon$  and hence

$$Q_{U;k}^{(t)} \leq 4B_0^8 \epsilon^2 v_1^2 \leq 4B_0^8 d^{-2} \log^4(T_0^* - 1);$$

We define event  $\mathcal{H}_{k;9;L} := (T_{B_0;k} > T_0^*) \setminus \mathcal{H}_{k;18;L}$ , then on the event  $\mathcal{H}_{k;9;L}$ , by applying Lemma 15 and 18, for all  $t \leq T_0^* \wedge T_{w;k}$  we have

$$\begin{aligned} U_k^{(t)} &\leq U_k^{(0)} + \sum_{j=4}^X \sum_{s=0}^{j-1} (v_1^{(s)})^2 - (v_k^{(s)})^2 U_k^{(s)} \\ &\leq T_0^* \leq 4B_0^8 d^{-2} \log^4(T_0^* - 1) + C_{16;L} \log^{5-2} \epsilon^{-1} B^4 d^{1-2} (T_0^*)^{1-2}. \end{aligned} \quad (\text{C.4})$$

Scaling condition (3.2) in Lemma 6 implies that

$$B^4 d^{1-2} (T_0^*)^{1-2} \log^4(T_0^* - 1) \leq \log^{5-2} \epsilon^{-1}. \quad (\text{C.5})$$

Together with (C.4), on the event  $\mathcal{H}_{k;9;L}$  we have

$$U_k^{(t)} \leq U_k^{(0)} + \sum_{j=4}^X \sum_{s=0}^{j-1} (v_1^{(s)})^2 - (v_k^{(s)})^2 U_k^{(s)} \leq C_{9;L} \log^{5-2} \epsilon^{-1} B^4 d^{1-2} (T_0^*)^{1-2};$$

where positive constant  $C_{9;L} = C_{16;L} + 4$ . Scaling condition (3.2) implies  $j \leq 3j((v_1^{(s)})^2 (v_k^{(s)})^2)^{1/2} [0; 1)$  for all  $s < T_{w;k}^0 \wedge T_{w;k}$ . From the reversed Gronwall Lemma 25, on the event  $H_{k;9;L}$  the following holds for all  $T_{w;k}^0 \wedge T_{w;k}$

$$U_k^{(t)} - U_k^{(0)} \leq \sum_{s=0}^{t-1} \left( 1 - j \leq 3j \left( (v_1^{(s)})^2 (v_k^{(s)})^2 \right)^{1/2} \right) \leq 2C_{9;L} \log^{5=2} \left( 1 - B^4 d^{1=2} (T_{w;k}^0)^{1=2} \right).$$

To obtain a lower bound on the probability of event  $H_{k;9;L}$ , we combine Lemmas 18, 19 and take union bound,

$$P(H_{k;9;L}) \geq 1 - 6 + 12 + \frac{5184}{\log^5 1}$$

■

### C.2. Proof of Lemma 10

Lemma 10 provides a quantitative characterization of the uniform initialization. As a probabilistic fact, the uniform distribution  $v^{(0)}$  in  $D_1$  is equal in distribution to  $k^{-1}$  with  $N(0; 1)$ . In addition, we have

$$\min_{2 \leq k \leq d} W_k^{(0)} = \min_{2 \leq k \leq d} \log \frac{(v_1^{(0)})^2}{(v_k^{(0)})^2} = \log(v_1^{(0)})^2 - \max_{2 \leq k \leq d} \log(v_k^{(0)})^2; \quad (C.6)$$

due to definition of  $W_k^{(t)}$  in (3.11) and the elementary inequality  $x - 1 \leq \log x$  for all  $x > 0$ . Intuitively, this means that  $\min_{2 \leq k \leq d} W_k^{(0)}$  is lower bounded by the spacing between the largest and second largest order statistics of i.i.d. logarithmic chi-squared distributions.

We provide an elementary probabilistic Lemma 20 to show the spacing on the right hand of (C.6)  $W_k^{(0)} = (\log^{-1} d)$  with high probability<sup>5</sup>.

Lemma 20 Let  $\{z_i^2\}_{i=1}^n$  be squares of i.i.d. standard normal variables, and denote their order statistics as  $z_{(1)}^2, \dots, z_{(n)}^2$ . Then for any  $\delta \in (0; 1=3)$ , when  $n \geq \frac{2}{\delta} \frac{1}{2e} \log^{-1} 1 + 1$ , we have with probability at least  $1 - \delta$

$$\log \frac{z_{(2)}^2}{z_{(n)}^2} \geq \log \frac{z_{(n-1)}^2}{z_{(n)}^2} \geq \frac{1}{8 \log^{-1} \log n}; \quad (C.7)$$

We can apply Lemma 20 and prove Lemma 10.

Proof [Proof of Lemma 10] Following assumptions in Theorem 7, we know that  $v^{(0)}$  has the same distribution as  $k^{-1}$  where  $v \sim N(0; 1)$ . Under scaling condition (3.2) in Lemma 6, Lemma 10 is straightforward if we apply Lemma 20 with  $n = d$  and use (C.6). ■

5. In retrospect, a similar lemma characterizing the lower bound of spacings was achieved using a different method in Bai et al. (2018).

C.3. Proof of Lemma 11

Let  $k \in [2; d]$  be any fixed coordinate. For each  $t \geq 1$ , we define random variable

$$Q_{W;k}^{(t)} = W_k^{(t)} - W_k^{(t-1)} + \text{sign}(v_4 - v_3) \cdot 2(v_1 - v_2)^3 v_1 v_k^3 (v_1 Y_k - v_k Y_1); \quad (\text{C.8})$$

Lemma 21 For each coordinate  $k \in [2; d]$  and any  $t \geq 1$ , on the event

$$H_{k;21;L}^{(t)} = \{ |v_1 - v_2| \leq B_0; |Y_1| \leq B_0; |Y_k| \leq B_0; v_1^2 < 3v_k^2; v_1^2 \leq \max_{i \in [2;d]} v_i^2 \}; \quad (\text{C.9})$$

under condition

$$12B^8 d^{-2} \log^4(T_{;1}^{(0)}) \leq 1; \quad (\text{C.10})$$

we have

$$|Q_{W;k}^{(t)}| \leq C_{21;L} B_0^8 d^{-2};$$

where  $C_{21;L}$  is a positive, absolute constant.

Lemma 22 Let Assumption 2 hold. For each coordinate  $k \in [2; d]$  and any  $t \geq 1$ , we have

$$\mathbb{E} \left[ \text{sign}(v_4 - v_3) \cdot 2(v_1 - v_2)^3 v_1 v_k^3 (v_1 Y_k - v_k Y_1) \mid \mathcal{F}_{t-1} \right] = 2 |v_4 - v_3| v_1^2 W_k; \quad (\text{C.11})$$

For each  $k \in [2; d]$  and  $t \geq 1$ , at the  $t$ -th iterate we let

$$f_k^{(t)} = \text{sign}(v_4 - v_3) \cdot 2(v_1 - v_2)^3 v_1 v_k^3 (v_1 Y_k - v_k Y_1) - 2 |v_4 - v_3| v_1^2 W_k; \quad (\text{C.12})$$

which, indexed by  $t$ , forms a sequence of martingale differences with respect to  $\mathcal{F}_t$ . Combining (C.8) and (C.11) together we have

$$W_k^{(t)} = \left( 1 + 2 |v_4 - v_3| v_1^2 \right)^t W_k^{(0)} + Q_{W;k}^{(t)} + f_k^{(t)}; \quad (\text{C.13})$$

By letting

$$P_t^0 = \prod_{s=0}^{t-1} \left( 1 + 2 |v_4 - v_3| v_1^2 \right)^{-1} \mathcal{F}_{t-1}; \quad (\text{C.14})$$

we conclude the following lemma.

Lemma 23 For each coordinate  $k \in [2; d]$  and any  $t \geq 1$ , the iteration  $W_k^{(t)}$  can be represented linearly as

$$P_t^0 W_k^{(t)} = W_k^{(0)} + \sum_{s=1}^t P_s^0 Q_{W;k}^{(s)} + \sum_{s=1}^t P_s^0 f_k^{(s)}; \quad (\text{C.15})$$

Lemma 24 Let Assumption 2 hold and initialization  $W^{(0)} \in D_{\text{cold}} \setminus D_{\text{mid};k}^c$ . Let  $\epsilon$  be a fixed positive. For each coordinate  $k \in [2; d]$ , there exists an event  $H_{k;24;L}$  satisfying

$$P(H_{k;24;L}) \geq 1 - \frac{6 + \frac{5184}{\log^5 1}}{\epsilon}$$

such that on event  $\mathcal{H}_{k;24;L}$  the following concentration result holds

$$\max_{T_{;0:5}^0 \wedge T_{c;k} \wedge T_1} \sum_{s=1}^t P_s^{\circ(s)} \mathbf{1}_k^{(s)} \leq C_{24;L} \log^{5=2} \frac{B^4}{j^4 3^{j=2}} d^{1=2}$$

where  $C_{24;L}$  is a positive, absolute constant.

Along with the tail probability Lemma 19, we are ready to present the proof of Lemma 11. Proof [Proof of Lemma 11] Recall the definition of truncation barrier  $\mathcal{B}_k$  in (C.1), stopping times  $T_{B_0;k}$  in (C.2),  $T_{c;k}$  in (3.14) and  $T_1$  in (3.15). We notice that (C.10) holds under scaling condition (3.2), and event  $(T_{;0:5}^0 \wedge T_{c;k} \wedge T_1) \setminus (T_{B_0;k} > T_{;0:5}^0) \cap \mathcal{H}_{k;21;L}^{(t)}$ . We define event  $\mathcal{H}_{k;11;L}$   $(T_{B_0;k} > T_{;0:5}^0) \setminus \mathcal{H}_{k;24;L}$ , then on event  $\mathcal{H}_{k;11;L}$ , by applying Lemma 21 for all  $T_{;0:5}^0 \wedge T_{c;k} \wedge T_1$  we obtain

$$v_1^2 \leq \frac{1}{d}; \quad |Q_{W;k}^{(t)}| \leq C_{21;L} B^8 d^2 \log^4(T_{;1}^0):$$

Scaling condition (3.2) also guarantees

$$\frac{j^4 3^j}{d} < 1; \quad \frac{B^4}{j^4 3^{j=2}} d^{1=2} \log^4(T_{;1}^0) \leq \log^{5=2} \frac{1}{d}; \quad (\text{C.16})$$

and hence by summation of geometric series, on event  $\mathcal{H}_{k;11;L}$  we have for all  $T_{;0:5}^0 \wedge T_{c;k} \wedge T_1$

$$\begin{aligned} & \sum_{s=1}^t P_s^{\circ(s)} Q_{W;k}^{(s)} \\ & \sum_{s=1}^t \left(1 + \frac{j^4 3^j}{d}\right)^s |Q_{W;k}^{(s)}| \leq \frac{1 + \frac{j^4 3^j}{d}}{1 - \frac{j^4 3^j}{d}} C_{21;L} B^8 d^2 \log^4(T_{;1}^0) \\ & = \frac{C_{21;L}}{2} \frac{B^8}{j^4 3^j} d^2 \log^4(T_{;1}^0) \leq \frac{C_{21;L}}{2} \log^{5=2} \frac{1}{d} \frac{B^4}{j^4 3^{j=2}} d^{1=2}. \end{aligned}$$

Combining with Lemmas 23 and 24, on event  $\mathcal{H}_{k;11;L}$  we know that for all  $T_{;0:5}^0 \wedge T_{c;k} \wedge T_1$  the following holds for constant  $C_{11;L} = C_{24;L} + C_{21;L} = 2$

$$W_k^{(t)} \leq \sum_{s=0}^t \left(1 + \frac{j^4 3^j}{d}\right)^s (v_1^{(s)})^2 \leq W_k^{(0)} \leq C_{11;L} \log^{5=2} \frac{1}{d} \frac{B^4}{j^4 3^{j=2}} d^{1=2}.$$

We verify the remaining claims in Lemma 11 by applying Lemma 19 with  $\frac{1}{d}$ , Lemma 24 and taking union bound

$$P(\mathcal{H}_{k;11;L}) \leq P(T_{B_0;k} > T_{;0:5}^0) + P(\mathcal{H}_{k;24;L}^c) \leq 15 + \frac{5184}{\log^5 \frac{1}{d}}.$$

■

C.4. Proof of Secondary Lemmas

Proof [Proof of Lemma 18]  $e_k^{(t)}$  earlier defined in (B.7) forms a martingale difference sequence with respect to  $\mathcal{F}_t$ . Recall (B.15) in proof of Lemma 16, we have derived the following Orlicz  $\psi_2$ -norm bound under Assumption 2

$$e_k^{(t)} \leq C_{28;L}^0 B^4 jv_{1j}^{-1} (1 + jv_{k=v_{1j}}):$$

Because initialization  $D_{mid;k}^{(0)} \leq 2$  on the event  $(t \leq T_{w;k})$ , where  $T_{w;k}$  is earlier defined in (3.9), we have  $jv_{1j}^{-1} d^{1=2}, jv_{k=v_{1j}}^{-1} = \bar{2}$ , and then

$$e_k^{(t)} \leq 2C_{28;L}^0 B^4 d^{1=2} :$$

Because  $1_{(t \leq T_{w;k})} \leq \mathcal{F}_t$ , we know that  $e_k^{(t)} 1_{(t \leq T_{w;k})}$  forms a martingale difference sequence with respect to  $\mathcal{F}_t$ , and  $k_{1=2} e_k^{(t)} k_{1=2} \leq 2C_{28;L}^0 B^4 d^{1=2}$ . Hence we can apply Theorem 29 with  $\psi = \psi_2$  and obtain for all  $\delta \in (0, e^{-1}]$ ,

$$\begin{aligned} P & \max_{1 \leq t \leq T_{w;k}^0} \sum_{s=1}^t e_k^{(s)} \leq C_{18;L} B^4 d^{1=2} (T_{w;k}^0)^{1=2} \log^{5=2} \frac{1}{\delta} \\ & = P \max_{1 \leq t \leq T_{w;k}^0} \sum_{s=1}^t e_k^{(s)} 1_{(s \leq T_{w;k})} \leq C_{18;L} B^4 d^{1=2} (T_{w;k}^0)^{1=2} \log^{5=2} \frac{1}{\delta} \\ & \leq 2 \cdot 3 + 6^4 \frac{64 T_{w;k}^0}{C_{18;L}^2 B^8 d^2 T_{w;k}^0 \log^5 \frac{1}{\delta}} \exp \left\{ -\frac{C_{18;L}^2 B^8 d^2 T_{w;k}^0 \log^5 \frac{1}{\delta}}{4C_{28;L}^0 B^8 d^2} \right\} \\ & = 6 + \frac{5184}{\log^5 \frac{1}{\delta}} ; \end{aligned}$$

where constant  $C_{18;L} \leq 8^p \bar{2} C_{28;L}^0$ . ■

Proof [Proof of Lemma 19] Applying Markov inequality and Lemma 12 gives

$$P(jv^{>Y} j > B_0) = P \left( \frac{jv^{>Y} j^2}{B^2} > \frac{B_0^2}{B^2} \right) \leq \exp \left( -\frac{B_0^2}{B^2} \right) \mathbb{E} \exp \left( \frac{jv^{>Y} j^2}{B^2} \right) \leq \frac{2}{T_{w;k}^0};$$

With the same procedure we also obtain

$$P(jY_1^{(t)} j > B_0) \leq \frac{2}{T_{w;k}^0}; \quad P(jY_k^{(t)} j > B_0) \leq \frac{2}{T_{w;k}^0};$$

Taking union bound,

$$\begin{aligned} P(T_{B_0;k} \leq T_{w;k}^0) & \leq \sum_{t=1}^{T_{w;k}^0} P(jv^{(t)} j > Y^{(t)} j > B_0) + P(jY_1^{(t)} j > B_0) + P(jY_k^{(t)} j > B_0) \\ & \leq 3T_{w;k}^0 \leq \frac{2}{T_{w;k}^0} \cdot 6(\delta + 1) ; \end{aligned}$$



due to  $T_{;1}^0 = T_{;1}^0 + 1$ , which comes from its definition in (3.1) and the elementary inequality  $\int_0^x e^{-(t+1)} dt \leq e^{-x}$  for all  $x \geq 0$ . ■

Proof [Proof of Lemma 20] Let  $F(x)$  be the cumulative distribution function of  $\chi^2(2)$ , where  $\chi^2(2)$  denote the chi-squared distribution, then

$$F(x) = 1 - e^{-x/2}; \quad F'(x) = \frac{1}{2} e^{-x/2}; \quad (C.17)$$

where  $F(x)$  and  $f(x)$  are cumulative distribution function and probability density function of standard Gaussian random variables. Let  $U_1, \dots, U_n$  be i.i.d. samples from  $U(0, 1)$ , whose largest and second largest order statistics  $U_{(n)}$  and  $U_{(n-1)}$  are denoted as  $U$  and  $M$  for notational simplicity. We also define  $M = F^{-1}(U)$ ,  $U = F^{-1}(M)$ . Since  $F(x)$  is concave for  $x \geq 0$ , under condition

$$M \geq 0; \quad \text{i.e. } U \geq F(0) = 0; \quad (C.18)$$

by concavity we have

$$U - M \leq F(M) - F(M) = F'(M)(M - M); \quad (C.19)$$

We denote the inverse functions of  $F$  by  $F^{-1}$ ,  $F^{-1}$ , then from (C.17) we have

$$F^{-1}(y) = 2 \log \left( \frac{1}{2} \frac{y+1}{2} \right); \quad (C.20)$$

(1) Taking the derivative for an inverse function gives

$$[F^{-1}]'(y) = \frac{1}{(F^{-1}(y))'} = \frac{1}{\frac{1}{2} \exp \left( -\frac{1}{2}(y+1) \right)};$$

Applying chain rule of derivative to (C.20), we have

$$[F^{-1}]'(y) = \frac{[F^{-1}]' \left( \frac{y+1}{2} \right)}{\frac{1}{2}} = \frac{\frac{1}{2} \exp \left( -\frac{1}{2} \frac{y+1}{2} \right)}{\frac{1}{2} \frac{y+1}{2}}; \quad (C.21)$$

(2) The following bound on tail probability of standard gaussian random variable is folklore:

$$\frac{x}{x^2 + 1} \leq \frac{1}{2} \exp \left( -\frac{x^2}{2} \right) \leq \frac{1}{x} \leq \frac{1}{2} \exp \left( -\frac{x^2}{2} \right); \quad (C.22)$$

We define

$$z_1 = \frac{x^2 + 1}{x} \exp \left( -\frac{x^2}{2} \right); \quad z_2 = x \exp \left( -\frac{x^2}{2} \right);$$

then for all  $x \geq 1$ , if  $z_1 \leq \frac{1}{2} \exp \left( -\frac{x^2}{2} \right)$ , i.e.  $x \geq \sqrt{2 \log z_1}$ , we have

$$\exp \left( -\frac{x^2}{2} \right) \leq \frac{z_1}{2x} \leq \frac{z_1}{2 \sqrt{2 \log z_1}};$$

which implies that

$$x \leq \frac{p}{2 \log z_1} \frac{\log \log z_1}{3 \log 2} \quad (C.23)$$

If  $z_2 = \exp(x^2 - 2)$ , we have

$$x \leq \frac{p}{2 \log z_2} \quad (C.24)$$

For all  $y \in [1 - \frac{p}{2e}; 1)$ , we can find the solution  $x$  of  $\frac{y+1}{2} = 1 - \frac{1}{z_1}$ : By applying (C.22) and (C.23) with this solution, we obtain

$$1 - \frac{y+1}{2} \leq x \leq \frac{p}{2 \log p - (1-y)} \frac{\log \log p - \frac{p}{2}}{3 \log 2} \quad (C.25)$$

We can find solution  $x$  to

$$\frac{y+1}{2} = 1 - \frac{1}{z_2}$$

By applying (C.22) and (C.24) with this solution, we obtain

$$1 - \frac{y+1}{2} \leq x \leq \frac{p}{2 \log p - (1-y)} \frac{p}{2} \quad (C.26)$$

Applying the lower bound (C.25) and upper bound (C.26) of  $\frac{y+1}{2}$  to (C.21), we have for all  $y \in [1 - \frac{p}{2e}; 1)$

$$[F^{-1}]^0(y) \leq \frac{1}{2 \log p - \frac{p}{2}} \frac{1}{8 \log p - \frac{p}{2}} \frac{1}{2(1-y) \log(1-y)^{-1}}$$

Replacing with  $U$  and  $F^{-1}(y)$  by  $M$  we have under condition

$$U \leq 1 - \frac{1}{2e} \quad (C.27)$$

we have

$$F^0(M) = \frac{1}{[F^{-1}]^0(U)} \frac{1}{2(1-U) \log(1-U)^{-1}} \quad (C.28)$$

- (3) It is standard from order statistics that  $U \sim \text{Beta}(1; n)$  and  $1-U \sim \text{Beta}(2; n-1)$ . Therefore, for any  $x \in (0; 1=3]$  and  $n \geq 2$ ,

$$P(U \leq \frac{x}{n}) = \int_0^{\frac{x}{n}} n(1-x)^{n-1} dx = 1 - (1 - \frac{x}{n})^n \geq 1 - \frac{1}{3} \quad (\log 3)$$

where we used elementary inequality  $(1-x)^x \geq 1-3x$  when  $x \leq \frac{1}{6}$ , and  $1 - 3^{-x} \geq (\log 3)x$  for all  $x \in (0; 1=3]$ . In addition, for any  $x \in (0; 1=3]$  and  $n \geq \frac{2}{2e \log^{-1} + 1}$  we have

$$\begin{aligned} P(1-U \leq \frac{2 \log^{-1}}{n-1}) &= \int_{\frac{2 \log^{-1}}{n-1}}^1 n(n-1)x(1-x)^{n-2} dx \\ &= 2 \log^{-1} \frac{n}{n-1} \frac{1}{n-1} + \frac{2 \log^{-1}}{n-1} \frac{1}{n-1} \\ &\geq 2 \log^{-1} + 1 \geq \frac{2 \log^{-1}}{n-1} \frac{1}{n-1} \geq (2 \log^{-1} + 1)^2 \geq \frac{2 \log 3 + 1}{3}; \end{aligned}$$

where we used elementary inequalities  $(1+x)^x \geq 1+e^{-x}$  for all  $x \geq 0$ , and  $(2 \log 1 + 1) \geq (2 \log 3 + 1) \geq 3$  for all  $2 \in (0; 1=3]$ . Taking union bound, we have

$$P \cup U \cup \frac{1}{n}; 1 \cup \frac{2 \log 1}{n-1} \geq 1-3: \tag{C.29}$$

(4) Notice that when  $2 \in (0; 1=3]$  and  $n \geq \frac{2^p}{2} e \log 1 + 1$ , conditions (C.18) and (C.27) hold automatically when (C.29) holds. Combining (C.19), (C.28) and (C.29), for  $2 \in (0; 1=3]$ ,  $n \geq \frac{2^p}{2} e \log 1 + 1$  we have with probability at least  $1-3$  that

$$M \cup M \cup \frac{U \cup U}{F^Q(M)} \cup \frac{U \cup U}{2(1-U) \log(1-U)} \geq \frac{(n-1)}{4n \log 1 \log \frac{n-1}{2 \log 1}} \geq \frac{1}{8 \log 1 \log n}:$$

Seeing that  $(M; M)$  is equidistributed as  $(\log \frac{2}{(n)}; \log \frac{2}{(n-1)})$ , for all  $2 \in (0; 1=3]$  and  $n \geq \frac{2^p}{2} e \log 1 + 1$  we have (C.7) holds with probability  $1-3$ . ■

Proof [Proof of Lemma 21] Recall that we define  $\epsilon_k = \text{sign}(v_k - Y_k)$ . For any fixed coordinate  $k \in [2; d]$ , we have

$$\begin{aligned} W_k^{(t)} - W_k^{(t-1)} &= \frac{v_1 + \epsilon_k (v > Y)^3 Y_1^2 - v_k + \epsilon_k (v > Y)^3 Y_k^2}{(v_k + \epsilon_k (v > Y)^3 Y_k)^2} - \frac{v_1^2 - v_k^2}{v_k^2} \\ &= \frac{2 \epsilon_k (v > Y)^3 v_1 v_k (v_k Y_1 - v_1 Y_k) + \frac{2}{3} (v > Y)^6 (v_k^2 Y_1^2 - v_1^2 Y_k^2)}{v_k^2 (v_k + \epsilon_k (v > Y)^3 Y_k)^2} \\ &= \epsilon_k \left( 1 + \frac{(v > Y)^3 Y_k}{v_k} \right)^2 \frac{v_k^4 - 2(v > Y)^3 v_1 v_k (v_k Y_1 - v_1 Y_k) + (v > Y)^6 (v_k^2 Y_1^2 - v_1^2 Y_k^2)}{v_k^2} : \end{aligned}$$

Combining with (C.8), this implies

$$\begin{aligned} Q_{W;k}^{(t)} &= 2 \epsilon_k \left( 1 + \frac{(v > Y)^3 Y_k}{v_k} \right)^2 \frac{1}{v_k} \left( (v > Y)^3 v_1 v_k^2 - Y_1 \frac{v_1}{v_k} Y_k \right. \\ &\quad \left. + \frac{2}{3} \left( 1 + \frac{(v > Y)^3 Y_k}{v_k} \right)^2 (v > Y)^6 v_k^2 - Y_1^2 \frac{v_1^2}{v_k^2} Y_k^2 \right) : \end{aligned}$$

Taylor series give for all  $|x| \leq \frac{1}{2}$  that

$$(1+x)^{-2} - 1 = \sum_{i=0}^{\infty} \binom{-2}{i} (-1)^i (i+2) x^i = \sum_{i=0}^{\infty} \binom{2}{i} x^i; \quad (1+x)^{-2} \leq 4:$$

On event  $H_{k;21;L}^{(t)}$ , since  $|v_1 - v_k| < \frac{1}{3}$  and  $|v_k - Y_k| \leq \frac{1}{3}$  ( $d \geq 3$ ), (C.10) implies  $\frac{(v > Y)^3 Y_k}{v_k} \leq \frac{1}{2}$ , and hence we have

$$|j Q_{W;k}^{(t)}| \leq 2 \left( 6 \frac{1}{3} B_0^4 d^{1=2} + (3 + 3 \frac{1}{3}) B_0^4 d^{1=2} + 2 \cdot 4 \cdot 12 B_0^8 d \right) \leq C_{21;L} B_0^8 d^2$$

where constant  $C_{21;L} = 156 + 36 \frac{1}{3}$ . ■

Proof [Proof of Lemma 22] Recall that  $W_k^{(t-1)} = (v_1^2 - v_k^2) = v_k^2$ . Under Assumption 2, using (B.13) we have

$$\begin{aligned} & E \sum_{i=1}^h \text{sign}(v_1 - v_k) 2v_1 v_k^3 (v_1 Y_k - v_k Y_1) F_{t-1}^i \\ &= 2 \text{sign}(v_1 - v_k) v_1 v_k^3 \sum_{i=1}^h (v_1 Y_k^3 + 3v_1 v_k \sum_{i=1}^h v_k v_1^3 - 3v_k v_1) \\ &= 2 \sum_{j=4}^3 v_1^2 W_k \end{aligned}$$

■

Proof [Proof of Lemma 23] From (C.13) and (C.14), we have

$$P_s^o W_k^{(s)} - P_{s-1}^o W_k^{(s-1)} = P_s^o Q_{W;k}^{(s)} + f_k^{(s)} \tag{C.30}$$

We iteratively apply (C.30) for  $s = 1, \dots, t$  and obtain (C.15). ■

Proof [Proof of Lemma 24]  $f_k^{(t)}$  defined in (C.12), forms a martingale difference sequence with respect to  $\mathcal{F}_{t-1}$ . Similar to techniques in proof of Lemma 16, we apply Lemma 27 three times, then use (E.1) and (B.12), in order to obtain the following bound on Orlicz norm

$$\begin{aligned} & 2(v_1 - v_k)^3 v_1 v_k^3 (v_k Y_1 - v_1 Y_k) \sum_{i=2}^h 2 \sum_{j=1}^i v_k^j v_1^j v_k^j \sum_{i=1}^2 (v_1 - v_k)^2 \sum_{i=1}^h (v_1 - v_k) Y_1 \frac{v_1}{v_k} Y_k \\ & 2 \sum_{j=1}^i v_k^j v_1^j v_k^j \sum_{i=2}^3 v_1 Y_1 \frac{v_1}{v_k} Y_k \sum_{i=2}^6 d^{1-2} v_1^3 \sum_{i=2}^3 (k Y_1 k + \sqrt{3} Y_k) \\ & 6(1 + \sqrt{3}) B^4 d^{1-2} \end{aligned}$$

where we use the fact that  $v_1 = v_k \sqrt{3}$  and  $v_k^2 = 1/(3d)$  on the event  $(t \leq T_{c;k} \wedge T_1)$ . Then by applying Lemma 28 we further derive

$$\begin{aligned} k f_k^{(t)} k_{1=2} &= 2(v_1 - v_k)^3 v_1 v_k^3 (v_k Y_1 - v_1 Y_k) E \sum_{i=1}^h 2(v_1 - v_k)^3 v_1 v_k^3 (v_k Y_1 - v_1 Y_k) \\ & C_{28;L}^0 2(v_1 - v_k)^3 v_1 v_k^3 (v_k Y_1 - v_1 Y_k) \sum_{i=2}^6 (1 + \sqrt{3}) C_{28;L}^0 B^4 d^{1-2} \end{aligned}$$

Because  $(t \leq T_{c;k} \wedge T_1) \subseteq \mathcal{F}_{t-1}$ , we know that  $f_k^{(t)} 1_{(t \leq T_{c;k} \wedge T_1)}$  forms a martingale difference sequence with respect to  $\mathcal{F}_{t-1}$ , and  $k f_k^{(t)} 1_{(t \leq T_{c;k} \wedge T_1)} k_{1=2} = k f_k^{(t)} k_{1=2}$ . Since  $v_1^2 = 1/d$  on the event  $(t \leq T_{c;k} \wedge T_1)$  and  $2 \sum_{j=4}^3 j = d < 1$  under scaling condition (3.2), by summation of geometric series, for all  $t \geq 1$  we have

$$\begin{aligned} & \sum_{s=1}^t (P_s^o)^2 f_k^{(s)} 1_{(t \leq T_{c;k} \wedge T_1)} \sum_{i=2}^6 1 + \frac{2 \sum_{j=4}^3 j}{d} \sum_{s=1}^{2s} 36(1 + \sqrt{3})^2 C_{28;L}^0 B^8 d^{-2} \\ & \frac{1 + \frac{2 \sum_{j=4}^3 j}{d}}{1 + \frac{2 \sum_{j=4}^3 j}{d}} \sum_{s=1}^{2s} 36(1 + \sqrt{3})^2 C_{28;L}^0 B^8 d^{-2} = 18(1 + \sqrt{3})^2 C_{28;L}^0 \frac{B^8}{\sum_{j=4}^3 j} d^2 \end{aligned}$$

With the bound given above, we apply Theorem 29 with  $\alpha = 2$  as

$$\begin{aligned}
 & P \max_{1 \leq t \leq T} \max_{0.5 \leq \tau \leq 1} \sum_{s=1}^t P_{s,k}^{(s)} C_{24;L} \log^{5-2} \frac{B^4}{j^4 3^{j-2}} d^{1-2} \\
 &= P \max_{1 \leq t \leq T} \sum_{s=1}^t P_{s,k}^{(s)} C_{24;L} \log^{5-2} \frac{B^4}{j^4 3^{j-2}} d^{1-2} \\
 & \leq 6 + 6^4 \frac{64 \cdot 18(1 + \frac{1}{3})^2 C_{28;L}^0 B^{8j-4} 3^{j-1} d^2}{C_{24;L}^2 \log^5 B^{8j-4} 3^{j-1} d^2} \\
 & \leq \exp \left\{ \frac{C_{24;L}^2 \log^5 B^{8j-4} 3^{j-1} d^2}{32 \cdot 18(1 + \frac{1}{3})^2 C_{28;L}^0 B^{8j-4} 3^{j-1} d^2} \right\}^{\frac{1}{5}} \\
 &= 6 + \frac{5184}{\log^5}
 \end{aligned}$$

where constant  $C_{24;L} = 24(1 + \frac{1}{3}) C_{28;L}^0$ . ■

### Appendix D. A Reversed Gronwall's Inequality

We present a discrete generalization of Gronwall's inequality, which is sharper than a straightforward application of Gronwall's inequality. Such sharp estimation plays a key role in our analysis. Although elementary, the lemma seems not recorded in relevant literature:

**Lemma 25** If for all  $t = 1, \dots, T$ ,  $\alpha(t) \in [0, 1)$  and if  $u(t)$  satisfies that for some positive constant

$$u(t) \leq u(0) + \sum_{0 \leq s < t} \alpha(s) u(s) \tag{D.1}$$

then for all  $t = 1, \dots, T$

$$u(t) \leq u(0) \prod_{0 \leq s < t} (1 - \alpha(s)) \leq \prod_{0 \leq s < t} (1 - \alpha(s)) \tag{D.2}$$

Note unlike analogous Gronwall-type results, here we pose no assumption on the sign of  $u(t)$ .  
**Proof** [Proof of Lemma 25] Define for all  $t = 1, \dots, T$

$$v(t) = \prod_{0 \leq s < t} (1 - \alpha(s))^{-1} u(0) + \sum_{0 \leq s < t} \alpha(s) u(s) \tag{D.3}$$

We have for  $s = 0, \dots, t-1$

$$\begin{aligned}
 v(s+1) - v(s) &= \sum_{0 \leq r \leq s+1} (1 - \alpha(r))^{-1} \alpha(s) u(s) \\
 &+ \sum_{0 \leq r \leq s} \alpha(u(0)) \sum_{0 \leq r \leq s+1} (r) u(r) A \sum_{0 \leq r \leq s+1} (1 - \alpha(r))^{-1} \sum_{0 \leq r \leq s+1} (1 - \alpha(r))^{-1} A \\
 &= \sum_{0 \leq r \leq s+1} (1 - \alpha(r))^{-1} \alpha(s) u(s) + \sum_{0 \leq r \leq s} \alpha(u(0)) \sum_{0 \leq r \leq s+1} (r) u(r) A \sum_{0 \leq r \leq s+1} (1 - \alpha(r))^{-1} \alpha(s) A \\
 &= \sum_{0 \leq r \leq s} \alpha(u(s)) u(0) + \sum_{0 \leq r \leq s+1} (r) u(r) A \sum_{0 \leq r \leq s+1} (1 - \alpha(r))^{-1} :
 \end{aligned}$$

Since  $\alpha(s) \in [0, 1)$ , and the product is nonnegative, the use of the lower side of (D.1) upper-estimates the difference  $v(s)$

$$v(s+1) - v(s) \leq \sum_{0 \leq r \leq s+1} (1 - \alpha(r))^{-1} \alpha(s) \tag{D.4}$$

Since  $v(0) = u(0)$ , telescoping the above inequality for  $s = 0, \dots, t-1$  gives

$$v(t) = v(0) + \sum_{0 \leq s < t} v(s+1) - v(s) \leq u(0) + \sum_{0 \leq s < t} \sum_{0 \leq r \leq s+1} (1 - \alpha(r))^{-1} \alpha(s)$$

and hence from the definition of  $v(t)$  in (D.3)

$$\begin{aligned}
 u(0) + \sum_{0 \leq s < t} \alpha(s) u(s) &= \sum_{0 \leq s < t} (1 - \alpha(s))^{-1} v(t) \\
 &= \sum_{0 \leq s < t} (1 - \alpha(s))^{-1} \alpha(u(0)) + \sum_{0 \leq s < t} \sum_{0 \leq r \leq s+1} (1 - \alpha(r))^{-1} A \\
 &= u(0) \sum_{0 \leq s < t} (1 - \alpha(s))^{-1} + \sum_{0 \leq s < t} \sum_{s+1 \leq r < t} (1 - \alpha(r)) :
 \end{aligned}$$

Taking the above result into the upper side of (D.1) gives

$$u(t) \leq u(0) + \sum_{0 \leq s < t} \alpha(s) u(s) + u(0) \sum_{0 \leq s < t} (1 - \alpha(s))^{-1} + \sum_{0 \leq s < t} \sum_{s+1 \leq r < t} (1 - \alpha(r)) ;$$

which further reduces to

$$u(t) \leq u(0) \sum_{0 \leq s < t} (1 - \alpha(s))^{-1} + \sum_{0 \leq s < t} \sum_{s+1 \leq r < t} (1 - \alpha(r)) :$$

That is, for all  $t = 0; 1; \dots; T$ ,

$$\begin{aligned}
 u(t) &= u(0) + \sum_{s=0}^{t-1} X(s) + \sum_{s=0}^{t-1} Y(s) \\
 &= u(0) + \sum_{s=0}^{t-1} X(s) + \sum_{s=0}^{t-1} Y(s) \\
 &= u(0) + \sum_{s=0}^{t-1} X(s) + \sum_{s=0}^{t-1} Y(s) \\
 &= u(0) + \sum_{s=0}^{t-1} X(s) + \sum_{s=0}^{t-1} Y(s)
 \end{aligned}$$

which proves the upper side of (D.2). For the lower side, applying the same inequality to the place of  $u$  gives the desired result. ■

### Appendix E. Orlicz $\psi_p$ -Norm

In this section, we recap Orlicz-norm and relevant properties. Many of the contributions can be found in standard texts [Wainwright \(2019\)](#); [van der Vaart and Wellner \(1996\)](#), and we derive its (old and new) conclusions for our use in online tensorial ICA analysis. We begin with its definition

**Definition 26 (Orlicz  $\psi_p$ -norm)** For a continuous, monotonically increasing and convex function  $\psi(x)$  defined for all  $x \geq 0$  satisfying  $\psi(0) = 0$  and  $\lim_{x \rightarrow \infty} \psi(x) = \infty$ , we define the Orlicz  $\psi_p$ -norm for a random variable  $X$  as

$$\|X\|_{\psi_p} = \inf \{K > 0 : E(\psi(|X|/K)) \leq 1\}.$$

When  $\psi(x)$  is monotonically increasing and convex for  $x \geq 0$ , the Orlicz  $\psi_p$ -norm satisfies triangle inequality, i.e. for any random variables  $X$  and  $Y$  we have  $\|X + Y\|_{\psi_p} \leq \|X\|_{\psi_p} + \|Y\|_{\psi_p}$ . The central case of interest in this paper is the Orlicz norm defined for a random variable  $X$  as

$$\|X\|_{\psi_2} = \inf \{K > 0 : E \exp \left( \frac{|X|^2}{K} \right) \leq 2\}.$$

In above we adopt the function  $\psi(x) = \exp(x^2) - 1$ . When  $p = 2$ ,  $\psi_p(x)$  is monotonically increasing and convex for  $x \geq 0$ , and consequently  $\psi_p$ -norm satisfies triangle inequality ([van der Vaart and Wellner, 1996](#))

$$\|X + Y\|_{\psi_p} \leq \|X\|_{\psi_p} + \|Y\|_{\psi_p} \tag{E.1}$$

We state a multiplicative property of the Orlicz  $\psi_p$ -norm which extends a standard result (e.g. Proposition D.3 in [Vu and Lei \(2013\)](#)):

**Lemma 27 (Multiplicative property)** Let  $X$  and  $Y$  be random variables then for some  $c_1$

$$\|XY\|_{\psi_p} \leq c_1 \|X\|_{\psi_p} \|Y\|_{\psi_p}.$$

Proof [Proof of Lemma 27] The inequality is trivial when  $\|X\| = 0$  or  $\|Y\| = 0$  since one of the variables is 0 a.s., or (ii) either  $X$  or  $Y$  has an infinite  $\ell_1$ -norm. Otherwise let  $\|X\| = \|Y\| = 1$ . Using  $\|X+Y\| \leq \frac{1}{4}(\|X\| + \|Y\|)^2$  and triangle inequality in (E.1) we have

$$\|X+Y\| \leq \frac{1}{4}(\|X\| + \|Y\|)^2 = \frac{1}{4}(\|X\| + \|Y\|)^2 = 1:$$

Multiplying both sides of the inequality by the positive  $\|X\| \|Y\|$  gives the desired result. ■

In the case of  $\alpha \in (0, 1)$ ,  $\phi_\alpha(x)$  is no longer convex; in fact it is only convex after a modification of the function in a neighborhood of  $x=0$ . Strictly speaking,  $\phi_\alpha$ -“norm” in this case is not a norm (defined in a Banach space) in the sense that the triangle inequality (E.1) does not hold for this set of  $\phi_\alpha$ 's. For the purpose of our applications, a generalized triangle inequality holds for the  $\alpha=2$  case in the following Lemma 28:

Lemma 28 For any random variables  $X, Y$  we have the following inequalities for Orlicz  $\psi_{1=2}$ -norm

$$\|X + Y\|_{\psi_{1=2}} \leq C_{28;L} (\|X\|_{\psi_{1=2}} + \|Y\|_{\psi_{1=2}}) \quad \text{and} \quad \|EX\|_{\psi_{1=2}} \leq C_{28;L} \|X\|_{\psi_{1=2}} \quad (E.2)$$

with  $C_{28;L} = 1.3937$ . In addition

$$\|X\|_{\psi_{1=2}} \leq \|EX\|_{\psi_{1=2}} \leq C_{28;L}^0 \|X\|_{\psi_{1=2}}; \quad (E.3)$$

where  $C_{28;L}^0 = 3.3359$

Proof [Proof of Lemma 28] We first rule out some simple cases. Analogously, we only need to prove (E.2) and (E.3) under the setting that  $\|X\|_{\psi_{1=2}} = 1$  and  $\|Y\|_{\psi_{1=2}} \in [0, 1]$ . Recall that when  $\alpha \in (0, 1)$ ,  $\phi_\alpha(x)$  does not satisfy convexity when  $x$  is around 0. Let the modified Orlicz-function be

$$\tilde{\phi}_\alpha(x) := \begin{cases} \exp(x) - 1 & x \geq x_0 \\ \frac{x}{x_0} (\exp(x_0) - 1) & x \in [0, x_0] \end{cases}$$

for some appropriate  $x_0 > 0$ , so as to make the function convex. Here  $x_0$  is chosen such that the tangent line of function  $\tilde{\phi}_\alpha$  at  $x_0$  passes through origin, i.e.

$$x_0^{-1} \exp(x_0) = \frac{\exp(x_0) - 1}{x_0}.$$

Simplifying it gives us a transcendental equation  $(\ln x_0 + 1) \exp(x_0) = 1$  which (admits no an analytic solution but) can be solved numerically. When  $\alpha = 2$ , we have  $x_{1=2} = 2.5396$ . Some numerical calculation yields

$$\|0\|_{\psi_{1=2}} \leq \tilde{\phi}_{1=2}(x) \leq 0.2666 \quad (E.4)$$

From (E.4) we immediately have

$$\|E \tilde{\phi}_{1=2}(jXj)\|_{\psi_{1=2}} \leq 1 \Rightarrow \|E \tilde{\phi}_{1=2}(jXj)\|_{\psi_{1=2}} \leq 1.2666 \quad \text{i.e.} \quad \|E \exp(jXj)^{1=2}\|_{\psi_{1=2}} \leq 2.2666$$



- (1) Let  $K_1, K_2$  denote the  $_{1=2}$  norms of  $X$  and  $Y$ , then  $E_{1=2}(jX=K_1j) = 1$  and  $E_{1=2}(jY=K_2j) = 1$ . Based on (E.4) we have  $E_{1=2}(j(X+Y)=(K_1+K_2)j) = 1$ . Applying triangle inequality (E.1) to Orlicz  $\psi_{1=2}$ -norm, we have

$$E_{1=2} \left( \frac{X+Y}{K_1+K_2} \right) = 1 \Rightarrow E_{1=2} \left( \frac{X+Y}{K_1+K_2} \right)^{1:2666} = 1:2666$$

where we applied (E.4). Now, applying Jensen's inequality to concave function  $f(x) = z^{\log_2 2:2666} x^2$  gives that, for constant  $C_{28;L} = (\log_2 2:2666)^2 = 1:3937$ ,

$$E_{1=2}(j(X+Y)=(C_{28;L}(K_1+K_2))j) = E \exp(j(X+Y)=(K_1+K_2)j^{1=2})^{\log_2 2:2666} = 1$$

$$E \exp(j(X+Y)=(K_1+K_2)j^{1=2})^{\log_2 2:2666} = 1 \quad 1;$$

which implies the first conclusion of (E.2):

$$\|kX + Yk_{1=2} \leq C_{28;L} (\|kXk_{1=2} + \|kYk_{1=2}):$$

- (2) Let  $K$  denote the  $_{1=2}$  norm of  $X$ , then  $E_{1=2}(jX=Kj) = 1$ . Based on (E.4) we have  $E_{1=2}(jEX=Kj) = 1$ . Because  $\psi_{1=2}$  is a convex function, we can apply Jensen's inequality as

$$\psi_{1=2}(jEX=Kj) \leq \psi_{1=2}(EjX=Kj) = E_{1=2}(jX=Kj) = 1$$

Combining with (E.4), it holds that  $E_{1=2}(jEX=Kj) = 1:2666$  which is equivalent to  $E_{1=2}(jEX=(C_{28;L}K)j) = 1$ , where  $C_{28;L} = (\log_2 2:2666)^2 = 1:3937$ . Noticing that  $E_{1=2}(EX=(CK)) = E_{1=2}(EX=(CK))$ , it holds that  $\|kEXk_{1=2} \leq C_{28;L} \|kXk_{1=2}$  which concludes (E.2).

- (3) To conclude (E.3), we have

$$\|kX - EXk_{1=2} \leq C_{28;L} (\|kXk_{1=2} + \|kEXk_{1=2}) \leq C_{28;L} (1 + C_{28;L}) \|kXk_{1=2} = C_{28;L}^0 \|kXk_{1=2}$$

where positive constant  $C_{28;L}^0 = C_{28;L} (1 + C_{28;L}) = 3:3359$  ■

## Appendix F. A Concentration Inequality for Martingales with Weak Exponential-type Tails

We prove a novel concentration inequality for 1-dimensional supermartingale difference sequence with finite  $\psi_2$ -norms that plays an important role in our analysis:

**Theorem 29** Let  $\beta \in (0; 1)$  be given. Assume that  $\{u_i : i = 1, \dots, N\}$  is a sequence of supermartingale differences with respect to  $\mathcal{F}_i$ , i.e.  $E[u_i | \mathcal{F}_{i-1}] = 0$ , and it satisfies  $\|ku_ik_{\psi_2} < 1$  for each  $i = 1, \dots, N$ . Then for an arbitrary  $N \geq 1$  and  $z > 0$ ,

$$P \left( \max_{1 \leq n \leq N} \sum_{i=1}^n u_i \geq z \right) \leq \frac{3}{3 + \frac{z^2}{64}} \exp \left( - \frac{z^2}{32 \sum_{i=1}^N \|ku_ik_{\psi_2}^2} \right) \leq \frac{9}{z^2} \sum_{i=1}^N \|ku_ik_{\psi_2}^2 \quad ; \quad (F.1)$$

Proof [Proof of Theorem 29] To prove Theorem 29, we will use a maxima version of the classical Azuma-Hoeffding's inequality proposed by Laib (1999) for bounded martingale differences, and then apply an argument of Lesigne and Wo (2001) and Fan et al. (2012) to truncate the tail and analyze the bounded and unbounded pieces separately.

- (1) First of all, for the sake of simplicity and with no loss of generality, throughout the following proof of Theorem 29 we shall pose the following extra condition

$$\sum_{i=1}^N ku_ik^2 = 1: \tag{F.2}$$

In other words, under the additional (F.2) condition proving (F.1) reduces to showing

$$P \left( \max_{i=1}^N u_i \geq z \right) \leq \frac{3}{3 + \frac{64}{z^2}} \exp \left( -\frac{z^2}{32} \right): \tag{F.3}$$

This can be made more clear from the following rescaling argument: one can put in the left of (F.3)  $u_i = \frac{u_i}{\sqrt{\sum_{i=1}^N ku_ik^2}}$  in the place of  $u_i$ , and  $z = \frac{z}{\sqrt{\sum_{i=1}^N ku_ik^2}}$  in the place of  $z$ , the left hand of (F.1) is just

$$P \left( \max_{i=1}^N \frac{u_i}{\sqrt{\sum_{i=1}^N ku_ik^2}} \geq \frac{z}{\sqrt{\sum_{i=1}^N ku_ik^2}} \right)$$

which, by (F.3), is upper-bounded by

$$\frac{3}{3 + \frac{64}{z^2}} \exp \left( -\frac{z^2}{32} \right) \leq \frac{9}{32 \sum_{i=1}^N ku_ik^2};$$

proving (F.1).

- (2) We apply a truncation argument used in Lesigne and Wo (2001) and later in Fan et al. (2012). Let  $M > 0$  be arbitrary, and we define

$$u_i^0 = u_i 1_{|u_i| \leq M}; \tag{F.4}$$

$$u_i^{00} = u_i 1_{|u_i| > M}; \tag{F.5}$$

$$T_n^0 = \sum_{i=1}^N u_i^0; \quad T_n^{00} = \sum_{i=1}^N u_i^{00}; \quad T_n^{000} = \sum_{i=1}^N E(u_i | \mathcal{F}_{i-1});$$

Since  $u_i$  is  $\mathcal{F}_i$ -measurable and  $u_i^0$  and  $u_i^{00}$  are two martingale difference sequences with respect to  $\mathcal{F}_i$ , and let  $T_n$  be defined as

$$T_n = \sum_{i=1}^N u_i \text{ and hence } T_n = T_n^0 + T_n^{00} + T_n^{000}. \tag{F.6}$$

Since  $u_i$  are supermartingale differences we have  $T_n^{000}$  is  $\mathcal{F}_{n-1}$ -measurable with  $T_0^{000} = 0$ ; a.s., and hence for any  $z > 0$ ,

$$\begin{aligned} P \max_{1 \leq n \leq N} T_n \geq z &= P \max_{1 \leq n \leq N} T_n^0 + T_n^{000} \geq z + P \max_{1 \leq n \leq N} T_n^{00} \geq z \\ &= P \max_{1 \leq n \leq N} T_n^0 \geq z + P \max_{1 \leq n \leq N} T_n^{00} \geq z \end{aligned} \quad (\text{F.7})$$

In the following, we analyze the tail bounds of  $T_n^0$  and  $T_n^{00}$  separately (Lesigne and Voln 2001; Fan et al., 2012).

(3) To obtain the first bound, we recap Laib's inequality as follows:

Lemma 30 (Laib, 1999) Let  $(w_i : 1 \leq i \leq N)$  be a real-valued martingale difference sequence with respect to some filtration  $\mathcal{F}_n$ , i.e.  $E[w_i | \mathcal{F}_{i-1}] = 0$ ; a.s., and the essential norm  $\|w_i\|_1$  is finite. Then for an arbitrary  $N \geq 1$  and  $z > 0$ ,

$$P \max_{1 \leq n \leq N} \sum_{i=1}^n w_i \geq z \leq \exp \left( -\frac{z^2}{2 \sum_{i=1}^N \|w_i\|_1^2} \right) \quad (\text{F.8})$$

(F.8) generalizes the folklore Azuma-Hoeffding's inequality, where the latter can be concluded from

$$\max_{1 \leq n \leq N} \sum_{i=1}^n w_i \leq \sum_{i=1}^N |w_i|$$

The proof of Lemma 30 is given in Laib (1999).

Recall our extra condition (F.2), then from the definition of  $u_i^0$  in (F.4) that  $\|u_i^0\|_1 \leq 2M \|u_i\|_1$ , the desired bound follows immediately from Laib's inequality in Lemma 30 by setting  $u_i^0$ :

$$P \max_{1 \leq n \leq N} T_n^0 \geq z = P \max_{1 \leq n \leq N} \sum_{i=1}^n u_i^0 \geq z \leq \exp \left( -\frac{z^2}{8M^2} \right) \quad (\text{F.9})$$

To obtain the tail bound of  $T_n^{00}$  we only need to show

$$E(u_i^{00})^2 \leq (6M^2 + 8B^2) \|u_i\|_1^2 \exp(-M \|u_i\|_1); \quad (\text{F.10})$$

where

$$B = \frac{3}{z}; \quad (\text{F.11})$$

from which, Doob's martingale inequality implies immediately that

$$P \max_{1 \leq n \leq N} T_n^{00} \geq z \leq \frac{1}{z^2} \sum_{i=1}^N E(u_i^{00})^2 \leq \frac{6M^2 + 8B^2}{z^2} \exp(-M \|u_i\|_1); \quad (\text{F.12})$$

To prove (F.10), first recall from the definition of  $u_i^{00}$  in (F.5) that

$$u_i^{00} = u_i \mathbf{1}_{\|u_i\|_1 > M \|u_i\|_1} \exp(-M \|u_i\|_1) + E[u_i \mathbf{1}_{\|u_i\|_1 > M \|u_i\|_1} | \mathcal{F}_{i-1}];$$

Recall from the property of conditional expectation that for any random variable  $W$  and a  $\sigma$ -algebra  $\mathcal{G} \subseteq \mathcal{F}$

$$\mathbb{E}[W - \mathbb{E}(W | \mathcal{G})]^2 = \mathbb{E}W^2 - \mathbb{E}[\mathbb{E}(W | \mathcal{G})]^2 \leq \mathbb{E}W^2 = \int_0^\infty 2y\mathbb{P}(|W| > y)dy$$

where the last equality is due to a simple application of Fubini's Theorem for nonnegative random variable  $|W|$ . Plugging in  $W = u_i 1_{\{|u_i| > \mathcal{M}\|u_i\| \}}$  and  $\mathcal{G} = \mathcal{F}_{i-1}$  we have

$$\begin{aligned} \mathbb{E}(u_i'')^2 &= \mathbb{E} \left[ u_i 1_{\{|u_i| > \mathcal{M}\|u_i\| \}} - \mathbb{E} \left( u_i 1_{\{|u_i| > \mathcal{M}\|u_i\| \}} \mid \mathcal{F}_{i-1} \right) \right]^2 \\ &\leq \int_0^\infty 2y\mathbb{P}(|u_i| 1_{\{|u_i| > \mathcal{M}\|u_i\| \}} > y)dy \\ &= \int_0^{\mathcal{M}\|u_i\|} 2ydy \cdot \mathbb{P}(|u_i| > \mathcal{M}\|u_i\|_\psi) + \int_{\mathcal{M}\|u_i\|}^\infty 2y\mathbb{P}(|u_i| > y)dy \quad (\text{F.13}) \\ &= \mathcal{M}^2 \|u_i\|_\psi^2 \mathbb{P}(|u_i| > \mathcal{M}\|u_i\|_\psi) + \int_{\mathcal{M}}^\infty 2t \|u_i\|_\psi \mathbb{P}(|u_i| > t \|u_i\|_\psi) \|u_i\|_\psi dt \\ &\leq 2\mathcal{M}^2 \|u_i\|_\psi^2 \exp\{-\mathcal{M}^\alpha\} + 4\|u_i\|_\psi^2 \int_{\mathcal{M}}^\infty t \exp\{-t^\alpha\} dt, \end{aligned}$$

where the last inequality is due to Markov's inequality that for all  $z > 0$

$$\mathbb{P}(|u_i|/\|u_i\|_\psi \geq z) \leq \exp\{-z^\alpha\} \mathbb{E} \exp\{|u_i|^\alpha/\|u_i\|_\psi^\alpha\} \leq 2 \exp\{-z^\alpha\}. \quad (\text{F.14})$$

It can be shown from basic calculus that the function  $g(t) = t^3 \exp\{-t^\alpha\}$  is decreasing in  $[\mathcal{B}, +\infty)$  and is increasing in  $[0, \mathcal{B}]$ , where  $\mathcal{B}$  was earlier defined in (F.11) (Fan et al., 2012). If  $\mathcal{M} \in [\mathcal{B}, \infty)$  we have

$$\begin{aligned} \int_{\mathcal{M}}^\infty t \exp\{-t^\alpha\} dt &= \int_{\mathcal{M}}^\infty t^{-2} t^3 \exp\{-t^\alpha\} dt \leq \int_{\mathcal{M}}^\infty t^{-2} dt \cdot \mathcal{M}^3 \exp\{-\mathcal{M}^\alpha\} \\ &= \mathcal{M}^{-1} \cdot \mathcal{M}^3 \exp\{-\mathcal{M}^\alpha\} = \mathcal{M}^2 \exp\{-\mathcal{M}^\alpha\}. \end{aligned} \quad (\text{F.15})$$

If  $\mathcal{M} \in (0, \mathcal{B})$ , we have by setting  $\mathcal{M}$  as  $\mathcal{B}$  in above

$$\begin{aligned} \int_{\mathcal{M}}^\infty t \exp\{-t^\alpha\} dt &= \int_{\mathcal{M}}^{\mathcal{B}} t \exp\{-t^\alpha\} dt + \int_{\mathcal{B}}^\infty t \exp\{-t^\alpha\} dt \\ &\leq \int_{\mathcal{M}}^{\mathcal{B}} dt \cdot \mathcal{B} \exp\{-\mathcal{M}^\alpha\} + \mathcal{B}^2 \exp\{-\mathcal{B}^\alpha\} \\ &\leq (\mathcal{B} - \mathcal{M})\mathcal{B} \exp\{-\mathcal{M}^\alpha\} + \mathcal{B}^2 \exp\{-\mathcal{M}^\alpha\} \leq 2\mathcal{B}^2 \exp\{-\mathcal{M}^\alpha\}. \end{aligned} \quad (\text{F.16})$$

Combining (F.13) with the two above displays (F.15) and (F.16) we obtain

$$\begin{aligned} \mathbb{E}(u_i'')^2 &\leq 2\mathcal{M}^2 \|u_i\|_\psi^2 \exp\{-\mathcal{M}^\alpha\} + 4\|u_i\|_\psi^2 \int_{\mathcal{M}}^\infty t \exp\{-t^\alpha\} dt \\ &\leq (6\mathcal{M}^2 + 8\mathcal{B}^2) \|u_i\|_\psi^2 \exp\{-\mathcal{M}^\alpha\}, \end{aligned}$$

completing the proof of (F.10) and hence (F.12).

- (4) Putting the pieces together: combining (F.7), (F.9) and (F.12) we obtain for an arbitrary  $u \in (0, \infty)$  that

$$\begin{aligned} \mathbb{P}\left(\max_{1 \leq n \leq N} T_n \geq 2z\right) &\leq \mathbb{P}\left(\max_{1 \leq n \leq N} T'_n \geq z\right) + \mathbb{P}\left(\max_{1 \leq n \leq N} T''_n \geq z\right) \\ &\leq \exp\left\{-\frac{z^2}{8\mathcal{M}^2}\right\} + \frac{6\mathcal{M}^2 + 8\mathcal{B}^2}{z^2} \exp\{-\mathcal{M}^\alpha\}. \end{aligned} \quad (\text{F.17})$$

We choose  $\mathcal{M}$  as, by making the exponents equal in above,

$$\mathcal{M} = \left(\frac{z^2}{8}\right)^{\frac{1}{\alpha+2}} \quad \text{such that} \quad \frac{z^2}{8\mathcal{M}^2} = \mathcal{M}^\alpha = \left(\frac{z^2}{8}\right)^{\frac{\alpha}{\alpha+2}}.$$

Plugging this  $\mathcal{M}$  back into (F.17) we obtain

$$\begin{aligned} \mathbb{P}\left(\max_{1 \leq n \leq N} T_n \geq 2z\right) &\leq \exp\left\{-\left(\frac{z^2}{8}\right)^{\frac{\alpha}{\alpha+2}}\right\} + \frac{6\mathcal{M}^2 + 8\mathcal{B}^2}{z^2} \exp\left\{-\left(\frac{z^2}{8}\right)^{\frac{\alpha}{\alpha+2}}\right\} \\ &\leq \left[1 + \left(\frac{1}{8}\right)^{\frac{2}{\alpha+2}} \frac{6}{z^{\frac{2}{\alpha+2}}} + \left(\frac{3}{\alpha}\right)^2 \frac{8}{z^2}\right] \exp\left\{-\left(\frac{z^2}{8}\right)^{\frac{\alpha}{\alpha+2}}\right\} \end{aligned} \quad (\text{F.18})$$

where we plugged in the expression of  $\mathcal{B}$  in (F.11). We can further simplify the square-bracket prefactor in the last line of (F.18) which can be tightly bounded by

$$\begin{aligned} 1 + \left(\frac{1}{8}\right)^{\frac{2}{\alpha+2}} \frac{6}{z^{\frac{2}{\alpha+2}}} + \left(\frac{3}{\alpha}\right)^2 \frac{8}{z^2} &\leq 1 + \frac{6 \cdot \frac{2}{\alpha+2}}{(8)^{\frac{2}{\alpha+2}}} + \frac{6 \cdot \frac{\alpha}{\alpha+2}}{(8)^{\frac{2}{\alpha+2}} z^2} + \left(\frac{3}{\alpha}\right)^2 \frac{8}{z^2} \\ &\leq 3 + \left(\frac{0.75 \cdot \frac{\alpha}{\alpha+2}}{(8)^{\frac{2}{\alpha+2}}} + \left(\frac{3}{\alpha}\right)^2\right) \frac{8}{z^2} \leq 3 + \left(0.75 + \left(\frac{3}{\alpha}\right)^2\right) \frac{8}{z^2} \leq 3 + \left(\frac{3}{\alpha}\right)^2 \frac{16}{z^2}. \end{aligned}$$

where we used an implication of Jensen's inequality: for  $\gamma = \alpha/(\alpha+2) \in (0, 1)$  one has  $x^\gamma \leq 1 - \gamma + \gamma x$  for all  $x \geq 0$  (where the equality holds for  $x = 1$ ), as well as a few elementary algebraic inequalities, including  $\gamma 8^{-\gamma} < 0.177$ ,  $(1 - \gamma)8^{-\gamma} < 1$ ,  $(3/\alpha)^{2/\alpha} > 0.78$  for all  $\alpha > 0$  and  $0 < \gamma = 2/(\alpha+2) < 1$ . Thus, (F.3) is concluded by noticing the relation (F.6) and setting  $z/2$  in the place of  $z$ , which hence proves Theorem 29 via the argument in (1) in our proof. ■

## Appendix G. Visualization and Numerical Experiments for Online Tensorial ICA

In this section, we provide visualization and conduct a group of experiments in support of our theoretical results. In §G.1 we illustrate the two-phase convergence behavior of our algorithm on two most commonly seen ICA problem distributions: Mixture Gaussian and Gaussian-Bernoulli. In §G.2 we empirically validate our sharp theoretical convergence rate (mainly in both  $d$  and  $T$ ) in our Corollary 8.

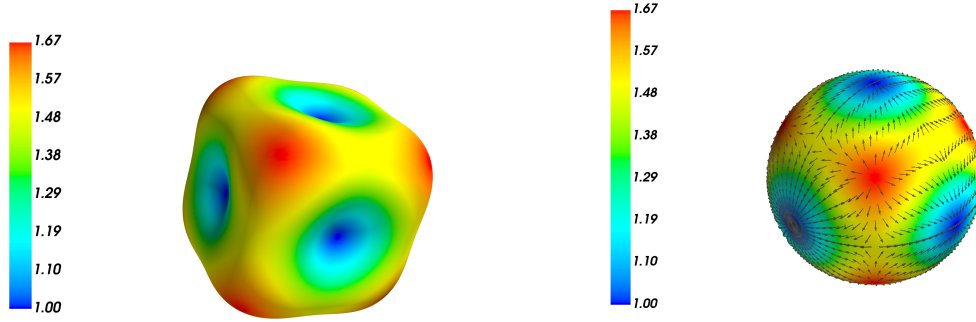


Figure 1: Heatmap visualization of the tensorial ICA objective function. Left: a graph plot where the radius denotes the corresponding function value (a global offset of 2 is added for illustration purposes). Right: a quiver plot with arrows denoting the negative gradient direction.

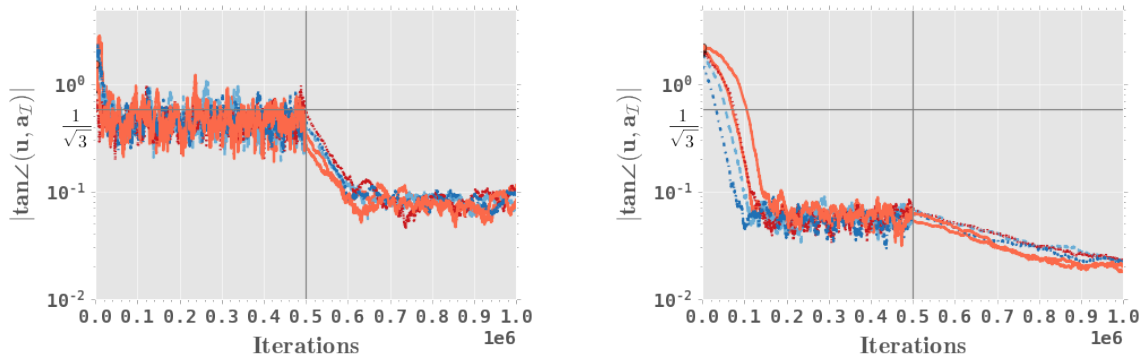


Figure 2: Convergence of our two-phase online tensorial ICA algorithm implemented under two settings. Left: each independent component being a mixture-Gaussian distribution with  $\mu_4 = 2.5$ . Right: each component being a Gaussian-Bernoulli distribution with  $\mu_4 = 6$ . Both settings are plotted in the semilog-scale, where the horizontal axis is the iteration number, and the vertical axis is the distance to the closest independent component in metric  $\min_{i \in [d]} |\tan \angle(\mathbf{u}^{(T)}, \mathbf{a}_i)|$ .

### G.1. Two-Phase Algorithm for Online Tensorial ICA

In this subsection, we present an initial experimental result that demonstrates the two-phase convergence behavior of our Algorithm 1. We arbitrarily set dimension as  $d = 20$  and total sample size as  $T = 1 \times 10^6$ , and the mixing matrix  $\mathbf{A}$  randomly drawn from the Haar measure over the  $d$ -dimensional orthogonal group. We conduct experiments on two separate instances of the one-dimensional distribution  $Z_i$  of independent component in ICA: mixture Gaussian with  $\mu_4 < 3$ , and Gaussian-Bernoulli with  $\mu_4 > 3$ , to validate the convergence theory of our two-phase algorithm.

- (a) **Mixture Gaussian** We adopt  $Z = \delta Y_1 + (1 - \delta)Y_2$  as the mixture Gaussian component of Gaussian variables  $Y_1, Y_2$  with

$$Y_1 \sim N\left(-\frac{1}{\sqrt{2}}, \frac{1}{2}\right), \quad Y_2 \sim N\left(\frac{1}{\sqrt{2}}, \frac{1}{2}\right), \quad \delta \sim \text{Bernoulli}(1/2) \quad \text{independently.}$$

It is easy to verify that the distribution of  $Z_i$  satisfies Assumption 2 with  $\mu_4 = 2.5$ .

- (b) **Gaussian-Bernoulli** where we adopt  $Z = \delta Y$  with

$$Y \sim N(0, 2), \quad \delta \sim \text{Bernoulli}(1/2) \quad \text{independently.}$$

It is straightforward to verify that  $\mu_4 = 6$  in this case.

The algorithm is initialized at  $\mathbf{u}_0$  uniformly drawn from the unit sphere  $\mathcal{D}_1$ , and we run our two-phase algorithm scheduled as in Corollary 8. That is, in the first half of the training when the number of iterates  $\leq T/2$  we choose stepsize as  $\frac{8d}{|\mu_4-3|T}$  and in the second half choose stepsize as  $\frac{9}{|\mu_4-3|T}$ , both omitting logarithmic factors. We measure the convergence by the tangent of the angle between  $\mathbf{u}^{(T)}$  and its closest independent component  $\mathbf{a}_i$ . The resulting 5 independent runs are shown in Figure 2. (The horizontal line represents  $y = \frac{1}{\sqrt{3}}$  which is the barrier of the warm region, and the vertical line is  $x = T/2$  which is the separation between two phases.) Our result in Figure 2 exemplifies consistency with the main theoretical result of our paper, especially the algorithm's two-phase demonstration. Regardless of the initialization all 10 independent runs share similar trajectories in terms of our measure: in Phase I the algorithm decays fast until oscillating into the warm region (the absolute tangent value below  $\frac{1}{\sqrt{3}}$ ), and in Phase II the algorithm continue to converge linearly until reaching the desired accuracy.

## G.2. Validating the Finite-Sample Convergence Rate

In this subsection, we validate our main convergence rate result, Corollary 8, via simulations for a range of values  $d$  and  $T$  satisfying the scaling condition

$$d \geq 2\sqrt{2\pi e} \log \epsilon^{-1} + 1 \quad \text{and} \quad \frac{d^4}{T} \leq C\epsilon^2, \quad (\text{G.1})$$

so with probability  $\geq 1 - O(\epsilon)$  the following  $\tilde{O}(\sqrt{d/T})$ -convergence rate result holds for all sufficiently large  $d$  and  $T$ :

$$|\tan \setminus(\mathbf{u}^{(T)}, \mathbf{a}_{\mathcal{I}})| \leq C\sqrt{\frac{d}{T}}, \quad (\text{G.2})$$

where  $C$  includes a polylogarithmic factor in  $d, T, \epsilon$ .

We continue to generate random data  $\mathbf{Z}$  as we did in the Gaussian-Bernoulli distributions case in §G.1, with parameter value of either  $d$  or  $T$  allowed to vary in each run while freezing the other. By ergodicity we may empirically estimate  $\mathbb{E}|\tan \setminus(\mathbf{u}^{(T)}, \mathbf{a}_{\mathcal{I}})|$  by averaging the last few iterates of  $|\tan \setminus(\mathbf{u}^{(T)}, \mathbf{a}_{\mathcal{I}})|$  — the final  $0.6T$  iterates are taken. Then we scatter-plot the empirical  $\mathbb{E}|\tan \setminus(\mathbf{u}^{(T)}, \mathbf{a}_{\mathcal{I}})|$  under a range of parameters in the figure, gauged by the theoretical rate of (G.2). We first run for  $T/2$  epochs with the stepsize propotional to  $8d/T$  and run for  $T/2$  epochs with the stepsize propotional to  $9/T$ , as shown in §G.1 and in Corollary 8, forgoing the logarithmic factor. The range of values one takes is  $d \in \{7, 12, 20, 33, 54\}$ ,  $T \in \{2 \times 10^2, 2 \times 10^3, 1 \times$

