This article was downloaded by: [2607:f470:a:7:a435:242b:d436:8d5e] On: 08 June 2022, At: 13:43

Publisher: Institute for Operations Research and the Management Sciences (INFORMS)

INFORMS is located in Maryland, USA



# **Operations Research**

Publication details, including instructions for authors and subscription information: <a href="http://pubsonline.informs.org">http://pubsonline.informs.org</a>

# Nonconvex Low-Rank Tensor Completion from Noisy Data

Changxiao Cai, Gen Li, H. Vincent Poor, Yuxin Chen

#### To cite this article:

Changxiao Cai, Gen Li, H. Vincent Poor, Yuxin Chen (2022) Nonconvex Low-Rank Tensor Completion from Noisy Data. Operations Research 70(2):1219-1237. <a href="https://doi.org/10.1287/opre.2021.2106">https://doi.org/10.1287/opre.2021.2106</a>

Full terms and conditions of use: <a href="https://pubsonline.informs.org/Publications/Librarians-Portal/PubsOnLine-Terms-and-Conditions">https://pubsonline.informs.org/Publications/Librarians-Portal/PubsOnLine-Terms-and-Conditions</a>

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact permissions@informs.org.

The Publisher does not warrant or guarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a guarantee, endorsement, or support of claims made of that product, publication, or service.

Copyright © 2021 The Author(s)

Please scroll down for article-it is on subsequent pages



With 12,500 members from nearly 90 countries, INFORMS is the largest international association of operations research (O.R.) and analytics professionals and students. INFORMS provides unique networking and learning opportunities for individual professionals, and organizations of all types and sizes, to better understand and use O.R. and analytics tools and methods to transform strategic visions and achieve better outcomes.

For more information on INFORMS, its publications, membership, or meetings visit <a href="http://www.informs.org">http://www.informs.org</a>

Vol. 70, No. 2, March-April 2022, pp. 1219-1237 ISSN 0030-364X (print), ISSN 1526-5463 (online)

#### **Methods**

# **Nonconvex Low-Rank Tensor Completion from Noisy Data**

### Changxiao Cai,<sup>a</sup> Gen Li,<sup>b</sup> H. Vincent Poor,<sup>a</sup> Yuxin Chen<sup>a</sup>

<sup>a</sup> Department of Electrical Engineering, Princeton University, Princeton, New Jersey 08540; <sup>b</sup> Department of Electronic Engineering, Tsinghua University, Beijing 100084, China

Contact: ccai@princeton.edu (CC); g-li16@mails.tsinghua.edu.cn (GL); poor@princeton.edu (HVP); yuxin.chen@princeton.edu, 

https://orcid.org/0000-0001-9256-5815 (YC)

Received: November 10, 2019
Revised: July 26, 2020
Accepted: November 1, 2020
Published Online in Articles in Advance:

Published Online in Articles in Advance June 3, 2021

.

Area of Review: Machine Learning and Data

Science

https://doi.org/10.1287/opre.2021.2106

Copyright: © 2021 The Author(s)

**Abstract:** We study a noisy tensor completion problem of broad practical interest, namely, the reconstruction of a low-rank tensor from highly incomplete and randomly corrupted observations of its entries. Whereas a variety of prior work has been dedicated to this problem, prior algorithms either are computationally too expensive for large-scale applications or come with suboptimal statistical guarantees. Focusing on "incoherent" and well-conditioned tensors of a constant canonical polyadic rank, we propose a two-stage nonconvex algorithm—(vanilla) gradient descent following a rough initialization—that achieves the best of both worlds. Specifically, the proposed nonconvex algorithm faithfully completes the tensor and retrieves all individual tensor factors within nearly linear time, while at the same time enjoying near-optimal statistical guarantees (i.e., minimal sample complexity and optimal estimation accuracy). The estimation errors are evenly spread out across all entries, thus achieving optimal  $\ell_{\infty}$  statistical accuracy. We also discuss how to extend our approach to accommodate asymmetric tensors. The insight conveyed through our analysis of nonconvex optimization might have implications for other tensor estimation problems.

Open Access Statement: This work is licensed under a Creative Commons Attribution 4.0 International License. You are free to copy, distribute, transmit and adapt this work, but you must attribute this work as "Operations Research. Copyright © 2021 The Author(s). https://doi.org/10.1287/opre.2021. 2106, used under a Creative Commons Attribution License: https://creativecommons.org/licenses/by/40/"

Funding: Y. Chen is supported in part by the Air Force Office of Scientific Research [Grant FA9550-19-1-0030], by the Office of Naval Research [Grant N00014-19-1-2120], by the Army Research Office [Grants W911NF-18-1-0303 and W911NF-20-1-0097], by the National Science Foundation (NSF) [Grants CCF-1907661, IIS-1900140, and IIS-2100158], and by a Princeton SEAS Innovation Award. H. V. Poor is supported in part by the NSF [Grant DMS-1736417]. C. Cai is supported in part by a Gordon Y. S. Wu Fellowship in Engineering. This work was done in part while Y. Chen was visiting the Kavli Institute for Theoretical Physics, supported in part by the NSF [Grant PHY-1748958]. Supplemental Material: The e-companion is available at https://doi.org/10.1287/opre.2021.2106.

Keywords: tensor completion • nonconvex optimization • gradient descent • spectral methods • entrywise statistical guarantees • minimaxity

# Introduction and Motivation Tensor Completion from Noisy Entries

Estimation of low-complexity models from highly incomplete observations is a fundamental task that spans a diverse array of science and engineering applications. Arguably one of the most extensively studied problems of this kind is matrix completion, where one wishes to recover a low-rank matrix given only partial entries (Davenport and Romberg 2016, Chen and Chi 2018). Moving beyond matrix-type data, a natural higher-order generalization is *low-rank tensor completion*, which aims to reconstruct a low-rank tensor when the vast majority of its entries are unseen. There is certainly no shortage of applications that motivate the investigation of tensor completion (e.g., personalized medicine (Soroushmehr and Najarian 2016,

Pawlowski 2019), medical imaging (Gandy et al. 2011, Semerci et al. 2014, Cheng et al. 2017), seismic data analysis (Kreimer et al. 2013, Ely et al. 2013), and multidimensional harmonic retrieval (Chen and Chi 2014, Ying et al. 2017)). One concrete example in operations research arises when learning the preference of individual customers for a collection of products on the basis of historical transactions (Farias and Li 2019, Mišić and Perakis 2020). Given the limited availability of transaction data (e.g., each customer might only have purchased very few products before), it is crucial to exploit multiway customer-product interactions (e.g., users' browsing and searching histories) in order to better predict the likelihood of a customer purchasing a new product. Clearly, the presence of missing data and the need of exploiting a multiway structure result in the

task of tensor completion. Additionally, tensor completion finds important applications in visual data inpainting (Liu et al. 2013, Li et al. 2017), where one wishes to reconstruct video data (or a sequence of images) from incomplete measurements. The video data consist of at least two spatial variables and one temporal variable, whose intrinsic connections are often modeled via certain low-complexity tensors.

For the sake of clarity, we phrase the problem formally before we proceed, focusing on a simple model that already captures the intrinsic difficulty of tensor completion in many aspects. Imagine that we are asked to estimate a symmetric order-3 tensor  $T^* \in \mathbb{R}^{d \times d \times d}$  from a small number of noisy entries

$$T_{j,k,l} = T_{i,k,l}^{\star} + E_{j,k,l}, \quad \forall (j,k,l) \in \Omega, \tag{1}$$

where  $T_{j,k,l}$  is the observed noisy entry at location (j,k,l),  $E_{j,k,l}$  stands for the associated noise, and  $\Omega \subseteq \{1,\cdots,d\}^3$  is a symmetric index subset to sample from. For notational simplicity, we set  $T = [T_{j,k,l}]_{1 \le j,k,l \le d}$  and  $E = [E_{j,k,l}]_{1 \le j,k,l \le d}$ , with  $T_{j,k,l} = E_{j,k,l} = 0$  for any  $(j,k,l) \notin \Omega$ . We adopt a random sampling model such that each index (j,k,l)  $(j \le k \le l)$  is included in  $\Omega$  independently with probability p. In addition, we know a priori that the unknown tensor  $T^* \in \mathbb{R}^{d \times d \times d}$  is a superposition of r rank-one tensors (often termed canonical polyadic [CP] decomposition if r is minimal):

$$T^{\star} = \sum_{i=1}^{r} u_{i}^{\star} \otimes u_{i}^{\star} \otimes u_{i}^{\star}, \text{ or more concisely,}$$

$$T^{\star} = \sum_{i=1}^{r} u_{i}^{\star \otimes 3}, \tag{2}$$

where each  $u_i^* \in \mathbb{R}^d$  represents one of the r low-rank tensor components/factors. Here and throughout, for any vectors  $a,b,c \in \mathbb{R}^d$ , the tensor  $a \otimes b \otimes c$  is a  $d \times d \times d$  array whose (j,k,l) th entry is given by  $a_jb_kc_l$ . The primary question is this: Can we hope to faithfully estimate  $T^*$ , as well as the individual tensor factors  $\{u_i^*\}_{1 \leq i \leq r}$ , from the partially revealed entries (1), assuming that r is reasonably small?

### 1.2. Computational and Statistical Challenges

Even though tensor completion conceptually resembles matrix completion in various ways, it is considerably more challenging than the matrix counterpart. This is perhaps not surprising, given that a plethora of natural tensor problems (e.g., computing the spectral norm, finding the best low-rank approximation) are all notoriously hard (Hillar and Lim 2013). As a notable example, whereas matrix completion is often efficiently solvable under nearly minimal sample complexity (Candès and Recht 2009, Gross 2011), all polynomial-time algorithms developed so far for tensor completion—even in the noise-free case—require a sample size at least exceeding the order of  $rd^{3/2}$ , which is substantially larger than the degrees of freedom (i.e., rd) underlying the model (2).

In fact, it is widely conjectured that there exists a large computational barrier away from the information-theoretic sampling limits (Barak and Moitra 2016).

With this fundamental gap in mind, the current paper focuses on the regime (in terms of the sample size) that enables reliable tensor completion in polynomial time. A variety of algorithms have been proposed that enjoy some sort of theoretical guarantees in (at least part of) this regime, including, but not limited to, spectral methods (Montanari and Sun 2018, Cai et al. 2021), sum-of-squares hierarchy (Barak and Moitra 2016, Potechin and Steurer 2017), nonconvex algorithms (Jain and Oh 2014, Xia and Yuan 2017), and also convex relaxation (based on proper unfolding) (Gandy et al. 2011, Romera-Paredes and Pontil 2013, Goldfarb and Qin 2014, Huang et al. 2015). Whereas these are all polynomial-time algorithms, most of the computational complexities supported by prior theory remain prohibitively high when dealing with large-scale tensor data a point that we shall elaborate on later. The only exception is the unfolding-based spectral method, which, however, fails to achieve exact recovery as the noise vanishes. This leads to the following critical question.

**Question 1.** *Is there any linear-time algorithm that is guaranteed to work for low-rank tensor completion?* 

Going beyond such computational concerns, one might naturally wonder whether it is also possible for a fast algorithm to achieve a nearly unimprovable statistical accuracy in the presence of noise. Toward this end, intriguing stability guarantees have been established for sum-of-squares hierarchy in the noisy settings (Barak and Moitra 2016), although this paradigm is computationally expensive for large-scale data. In a recent work, Xia et al. (2017) came up with a two-stage algorithm (i.e., a spectral method followed by tensor power iterations) for noisy tensor completion. Its estimation accuracy, however, falls short of achieving exact recovery in the absence of noise. This gives rise to another question of fundamental importance.

**Question 2.** Can we achieve near-optimal statistical accuracy without compromising computational efficiency?

In this paper, we aim to address these two questions by developing a nonconvex algorithm that achieves optimal computational efficiency and statistical accuracy all at once.

# 2. Algorithm and Main Results

## 2.1. A Two-Stage Nonconvex Algorithm

To address the aforementioned challenges, a first impulse is to resort to the following least-squares problem:

$$\underset{u_1, \dots, u_r \in \mathbb{R}^d}{\text{minimize}} \quad \sum_{j,k,l \in \Omega} \left( \left[ \sum_{i=1}^r u_i^{\otimes 3} \right]_{j,k,l} - T_{j,k,l} \right)^2, \tag{3}$$

or, more concisely (up to proper rescaling),

$$\underset{\boldsymbol{U} \in \mathbb{R}^{d \times r}}{\text{minimize}} \quad f(\boldsymbol{U}) := \frac{1}{6p} \left\| \mathcal{P}_{\Omega} \left( \sum_{i=1}^{r} \boldsymbol{u}_{i}^{\otimes 3} - T \right) \right\|_{F}^{2}, \tag{4}$$

if we take  $U := [u_1, \ldots, u_r] \in \mathbb{R}^{d \times r}$ . Here, we denote by  $\mathcal{P}_{\Omega}(T)$  the orthogonal projection of any tensor T onto the subspace of tensors that vanish outside of the index set  $\Omega$ . This optimization problem, however, is highly nonconvex (which involves minimizing a degree-6 polynomial), thus resulting in computational intractability in general.

Fortunately, not all nonconvex problems are as daunting to solve as they may seem. For example, recent years have seen a flurry of activity in low-rank matrix factorization via nonconvex optimization, which provably achieves optimal statistical accuracy and computational efficiency at once (see Chi et al. 2019 for an overview of recent advances). Motivated by this strand of work, we propose to solve (4) via a two-stage nonconvex paradigm, which we present in reverse order. The whole procedure is summarized in Algorithms 1–3.

**2.1.1. Gradient Descent (GD).** Arguably one of the simplest optimization algorithms is gradient descent, which adopts a gradient update rule

$$\mathbf{U}^{t+1} = \mathbf{U}^t - \eta_t \nabla f(\mathbf{U}^t), \quad t = 0, 1, \dots,$$
 (5)

where  $\eta_t$  is the learning rate or the step size and  $\boldsymbol{U}^t \in \mathbb{R}^{d \times r}$  is the estimate in the t th iteration. The main computational burden in each iteration lies in gradient evaluation, which, in this case, can be performed in time proportional to that taken to read the data.

Despite the simplicity of this algorithm, two critical issues stand out and might significantly affect its efficiency, which we shall bear in mind throughout the algorithmic and theoretical development.

i. Local stationary points and initialization. As is well known, GD is guaranteed to find an approximate local stationary point, provided that the learning rates do not exceed the inverse Lipschitz constant of the gradient (Bubeck 2015). There exist, however, local stationary points (e.g., saddle points or spurious local minima) that might fall short of the desired statistical properties. This requires us to properly avoid such undesired points, while retaining computational efficiency. To address this issue, one strategy is to first identify a rough initial guess within a local region surrounding the global solution (which often helps rule out bad local minima), in order to guarantee proper convergence of subsequent optimization procedures (Jain and Oh 2014, Li and Tang 2017). As a side remark, although careful initialization might not be crucial for several matrix recovery cases (Gilboa et al. 2018, Chen et al. 2019c, Tan and Vershynin 2019), it does seem to be critical in various tensor problems (Richard and Montanari 2014). We shall elucidate this point in Section EC.1 in the e-companion.

ii. Learning rates and regularization. Learning rates play a pivotal role in determining the convergence properties of GD. The challenge, however, is that the loss function (4) is overall not sufficiently smooth (i.e., its gradient often has an exceedingly large Lipschitz constant), and hence generic optimization theory recommends a pessimistically slow update rule (i.e., an extremely small learning rate) so as to guard against overshooting. This, however, slows down the algorithm significantly, thus destroying the main computational advantage of GD (i.e., low per-iteration cost). With this issue in mind, prior literature suggests carefully designed regularization steps (e.g., proper projection, regularized loss functions) in order to improve the geometry of the optimization landscape (Xia and Yuan 2017). In contrast, we argue that one is allowed to take a constant learning rate which is as aggressive as it can possibly be—even without enforcing any regularization procedures.

**Algorithm 1** (Gradient Descent for Nonconvex Tensor Completion)

1: Generate an initial estimate  $U^0 \in \mathbb{R}^{d \times r}$  via Algorithm 2.

2: **for** 
$$t = 0, 1, ..., t_0 - 1$$
 **do**

3: 
$$\boldsymbol{U}^{t+1} = \boldsymbol{U}^t - \eta_t \nabla f(\boldsymbol{U}^t) = \boldsymbol{U}^t - \frac{\eta_t}{p} \mathcal{P}_{\Omega}(\sum_{i=1}^r (\boldsymbol{u}_i^t)^{\otimes 3} - T)$$
  
 $\times_1^{\text{seq}} \boldsymbol{U}^t \times_2^{\text{seq}} \boldsymbol{U}^t$ , where  $\times_1^{\text{seq}}$  and  $\times_2^{\text{seq}}$  are defined in Section 2.4

**2.1.2. Initialization.** Motivated by the aforementioned issue (i), we develop a procedure that guarantees a reasonable initial estimate. In a nutshell, the proposed procedure consists of two steps:

a. estimate the subspace spanned by the r low-rank tensor factors  $\{u_i^{\star}\}_{1 \leq i \leq r}$  via a spectral method;

b. disentangle individual low-rank tensor factors from this subspace estimate.

As we shall see momentarily, the total computational complexity of the proposed initialization is  $O(pd^3)$  when r=O(1),  $\kappa=O(1)$ , and  $p\geq 1/d^2$  (where  $\kappa$  is a sort of "condition number" defined later), which is a linear-time algorithm. Note, however, that these two steps in the initialization procedure are relatively more complicated to describe. To improve the flow of the current paper, we postpone the details to Section 3. The readers can catch a glimpse of these procedures in Algorithms 2–3.

**Algorithm 2** (Spectral Initialization for Nonconvex Tensor Completion)

1: Let  $U \wedge U^{\top}$  be the rank-r eigen-decomposition of

$$B := \mathcal{P}_{\mathsf{off-diag}}(AA^{\top}), \tag{6}$$

where  $A = \text{unfold}(p^{-1}T)$  is the mode-1 matricization of  $p^{-1}T$  and  $\mathcal{P}_{\text{off-diag}}(\mathbf{Z})$  extracts out the off-diagonal entries of  $\mathbf{Z}$ .

2: **Output:** an initial estimate  $U^0 \in \mathbb{R}^{d \times r}$  on the basis of  $U \in \mathbb{R}^{d \times r}$  using Algorithm 3.

**Algorithm 3** (Retrieval of Low-Rank Tensor Factors from a Given Subspace Estimate)

1: **Input:** number of restarts L, pruning threshold  $\epsilon_{th}$ , subspace estimate  $U \in \mathbb{R}^{d \times r}$  given by Algorithm 2.

2: **for**  $\tau = 1, ..., L$  **do** 

3: Generate an independent Gaussian vector  $g^{\tau} \sim$  $\mathcal{N}(0, \mathbf{I}_d)$ .

4:  $(v^{\tau}, \lambda_{\tau}, \text{spec-gap}_{\tau}) \leftarrow \text{Retrieve-one-tensor-factor}$  $(T,p,U,g^{\tau}).$ 

 $\lambda_r$ )}  $\leftarrow \mathsf{PRUNE}(\{(\nu^\tau, \lambda_\tau, \mathsf{spec\text{-}gap}_\tau)\}_{\tau=1}^L, \epsilon_{\mathsf{th}}).$ 

6: **Output:** initial estimate  $\boldsymbol{U}^0 = [\lambda_1^{1/3} \boldsymbol{w}^1, \dots, \lambda_r^{1/3} \boldsymbol{w}^r].$ 

1: function Retrieve-one-tensor-factor(T, p, U, g)

2: Compute

$$\theta = UU^{\mathsf{T}}g =: \mathcal{P}_{U}(g), \tag{7a}$$

$$M = p^{-1}T \times_3 \theta, \tag{7b}$$

where  $\times_3$  is defined in Section 2.4.

3: Let  $\nu$  be the leading singular vector of M obeying  $\langle T, v^{\otimes 3} \rangle \ge 0$ , and set  $\lambda = \langle p^{-1}T, v^{\otimes 3} \rangle$ .

4: **return**  $(\boldsymbol{\nu}, \lambda, \sigma_1(\boldsymbol{M}) - \sigma_2(\boldsymbol{M}))$ .

 $\begin{array}{l} 1{:}\ \mathbf{function}\ \mathrm{Prune}(\{(\boldsymbol{\nu}^{\tau},\!\lambda_{\tau},\!\mathrm{spec-gap}_{\tau})\}_{\tau=1}^{L},\boldsymbol{\epsilon}_{\mathsf{th}}) \\ 2{:}\ \mathrm{Set}\ \boldsymbol{\Theta} = \{(\boldsymbol{\nu}^{\tau},\!\lambda_{\tau},\!\mathrm{spec-gap}_{\tau})\}_{\tau=1}^{L}. \end{array}$ 

3: **for** i = 1, ..., r **do** 

4: Choose ( $v^{\tau}$ ,  $\lambda_{\tau}$ , spec–gap<sub> $\tau$ </sub>) from Θ with the largest spec-gap<sub> $\tau$ </sub>; set  $w^i = v^{\tau}$  and  $\lambda_i = \lambda_{\tau}$ .

5: Update  $\Theta \leftarrow \Theta \setminus \{(v^{\tau}, \lambda_{\tau}, \text{spec-gap}_{\tau}) \in \Theta : |\langle v^{\tau}, w^{i} \rangle| >$  $1 - \epsilon_{\mathsf{th}}$ .

6: **return**  $\{(w^1, \lambda_1), \dots, (w^r, \lambda_r)\}.$ 

### 2.2. Main Results

Encouragingly, the proposed nonconvex algorithm provably achieves the best of both worlds—in terms of statistical accuracy and computational efficiency for a class of low-rank, well-conditioned, and "incoherent" problem instances. This subsection summarizes our main findings.

Before continuing, we note that one cannot hope to recover an arbitrary tensor from highly subsampled and arbitrarily corrupted entries. In order to enable provably valid recovery, the present paper focuses on a tractable model by imposing the following assumptions.

**Definition 1** (Incoherence and Well-Conditionedness). Define the incoherence parameters and the condition number of  $T^*$  as follows:

$$\mu_0 := \frac{d^3 \|T^*\|_{\infty}^2}{\|T^*\|_{\mathbb{F}}^2},\tag{8a}$$

$$\mu_1 := \frac{d||\boldsymbol{u}_i^{\star}||_{\infty}^2}{||\boldsymbol{u}_i^{\star}||_2^2},\tag{8b}$$

$$\mu_2 := \frac{d\langle u_i^*, u_j^* \rangle^2}{\|u_i^*\|_2^2 \|u_i^*\|_2^2},\tag{8c}$$

$$\kappa := \frac{\max_i \|\boldsymbol{u}_i^*\|_2}{\min_i \|\boldsymbol{u}_i^*\|_2}.$$
 (8d)

**Remark 1.** Here,  $\mu_0$ ,  $\mu_1$ , and  $\mu_2$  are termed the *incoher*ence parameters. Definitions (8a)–(8c) can be viewed as some sort of incoherence conditions for the tensor. For instance, when  $\mu_0, \mu_1$ , and  $\mu_2$  are small, these conditions say that (1) the energy of tensor  $T^*$  is (nearly) evenly spread across all entries; (2) each factor  $u_i^*$  is delocalized; (3) the factors  $\{u_i^*\}$  are nearly orthogonal to each other. Definition (8d) is concerned with the "well-conditionedness" of the tensor, meaning that each rank-1 component is of roughly the same size. In particular, we note that an assumption on pairwise correlation (i.e., a constraint on  $\mu_2$ ) is often assumed in the literature of tensor decomposition/factorization (e.g., Anandkumar et al. 2014b, Sun et al. 2017, Hao et al. 2020).

For notational simplicity, we shall set

$$\mu := \max\{\mu_0, \mu_1, \mu_2\}. \tag{9}$$

Note that our theory allows  $\mu$  to grow with the problem dimension d (in fact,  $\mu$  can be as large as d/polylog(d)).

**Assumption 1** (Random Noise). Suppose that E is a symmetric random tensor, where  $\{E_{j,k,l}\}_{1 \leq j \leq k \leq l \leq d}$  (see (1)) are independently generated sub-Gaussian random variables with mean zero and variance  $Var(E_{i,k,l}) \leq \sigma^2$ .

In addition, recognizing that there is a global permutational ambiguity issue (i.e., one cannot distinguish  $u_1^{\star}, \dots, u_r^{\star}$  from an arbitrary permutation of them), we introduce the following loss metrics to account for this ambiguity:

$$\mathsf{dist}_{\mathsf{F}}(\boldsymbol{U},\boldsymbol{U}^{\star}) := \min_{\boldsymbol{\Pi} \in \mathsf{perm}_r} \|\boldsymbol{U}\boldsymbol{\Pi} - \boldsymbol{U}^{\star}\|_{\mathsf{F}}, \tag{10a}$$

$$\mathsf{dist}_{\infty}(\boldsymbol{U}, \boldsymbol{U}^{\star}) := \min_{\Pi \in \mathsf{perm}_r} \|\boldsymbol{U}\Pi - \boldsymbol{U}^{\star}\|_{\infty}, \tag{10b}$$

$$\mathsf{dist}_{2,\infty}(\boldsymbol{U},\boldsymbol{U}^{\star}) := \min_{\boldsymbol{\Pi} \in \mathsf{perm}_r} \|\boldsymbol{U}\boldsymbol{\Pi} - \boldsymbol{U}^{\star}\|_{2,\infty}, \tag{10c}$$

where  $perm_r$  stands for the set of  $r \times r$  permutation matrices. For notational simplicity, we also take

$$\lambda_{\min}^{\star} := \min_{1 \le i \le r} \| \mathbf{u}_{i}^{\star} \|_{2}^{3} \quad \text{and} \quad \lambda_{\max}^{\star} := \max_{1 \le i \le r} \| \mathbf{u}_{i}^{\star} \|_{2}^{3}.$$
 (11)

With these notations in place, we are ready to present our main results. For simplicity of presentation, we shall start with the setting where  $r, \mu, \kappa \approx 1$ .

**Theorem 1.** Fix an arbitrary small constant  $\delta > 0$ . Suppose that  $r, \kappa, \mu = O(1)$ ,

$$p \ge c_0 \frac{\log^4 d}{d^{3/2}}, \qquad \frac{\sigma}{\lambda_{\min}^{\star}} \le c_1 \frac{\sqrt{p}}{d^{3/4} \log^2 d},$$

$$L = c_2 \quad \text{and} \quad \epsilon_{\text{th}} = c_3 \left( \frac{\log d}{d\sqrt{p}} + \frac{\sigma}{\lambda_{\min}^{\star}} \sqrt{\frac{d \log^2 d}{p}} + \sqrt{\frac{\log d}{d}} \right)$$

for some sufficiently large constants  $c_0$ ,  $c_2 > 0$  and some suffi*ciently small constants*  $c_1, c_3 > 0$ . The learning rate  $\eta_t \equiv \eta$  is taken to be a constant obeying  $0 < \eta \le \lambda_{\min}^{\star 4/3}/(32\lambda_{\max}^{\star 8/3})$ . Then with probability at least  $1 - \delta$ ,

$$\mathsf{dist}_{\mathsf{F}}(\boldsymbol{U}^{t}, \boldsymbol{U}^{\star}) \leq \left(C_{1} \rho^{t} + C_{2} \frac{\sigma}{\lambda_{\min}^{\star}} \sqrt{\frac{d \log d}{p}}\right) \|\boldsymbol{U}^{\star}\|_{\mathsf{F}}, \quad (12a)$$

 $\operatorname{dist}_{\infty}(\boldsymbol{U}^{t},\boldsymbol{U}^{\star}) \leq \operatorname{dist}_{2,\infty}(\boldsymbol{U}^{t},\boldsymbol{U}^{\star})$ 

$$\leq \left(C_3 \rho^t + C_4 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}}\right) \|\boldsymbol{U}^*\|_{2,\infty} \quad (12b)$$

hold simultaneously for all  $0 \le t \le t_0 = d^5$ . Here,  $0 < C_1, C_3, \rho < 1$  and  $C_2, C_4 > 0$  are some absolute constants.

**Remark 2.** The theorem holds unchanged if  $d^5$  is replaced by  $d^c$  for an arbitrarily large constant c > 0.

**Remark 3.** The upper bound  $t_0$  on the iteration count arises from the leave-one-out analysis when handling noisy observations. In short, the leave-one-out argument can only provide high-probability bounds for each iteration, thus requiring an upper bound on the iteration count if we desire a uniform bound across iterations. Note that in the noiseless case, our results and analysis hold for an arbitrarily large number of iterations.

As an immediate consequence of Theorem 1, we obtain appealing  $\ell_\infty$  statistical guarantees for estimating tensor entries, which are previously rarely available (see Table 1). Specifically, let our tensor estimate in the t th iteration be

$$T^{t} := \sum_{i=1}^{r} \boldsymbol{u}_{i}^{t} \otimes \boldsymbol{u}_{i}^{t} \otimes \boldsymbol{u}_{i}^{t}, \quad \text{where } \boldsymbol{U}^{t} = [\boldsymbol{u}_{1}^{t}, \dots, \boldsymbol{u}_{r}^{t}] \in \mathbb{R}^{d \times r}.$$
(13)

Then our result is the following.

**Corollary 1.** Fix an arbitrarily small constant  $\delta > 0$ . Instate the assumptions of Theorem 1. Then with probability at least  $1 - \delta$ ,

$$\| \mathbf{T}^t - \mathbf{T}^* \|_{\mathrm{F}} \lesssim \left( C_1 \rho^t + C_2 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{\rho}} \right) \| \mathbf{T}^* \|_{\mathrm{F}}, \quad (14a)$$

$$\| \mathbf{T}^t - \mathbf{T}^* \|_{\infty} \lesssim \left( C_3 \rho^t + C_4 \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{d \log d}{p}} \right) \| \mathbf{T}^* \|_{\infty} \quad (14b)$$

hold simultaneously for all  $0 \le t \le t_0 = d^5$ . Here,  $0 < C_1$ ,  $C_3$ ,  $\rho < 1$  and  $C_2$ ,  $C_4 > 0$  are some absolute constants.

Several important implications are provided as follows. The following discussion assumes that  $\lambda^{\star}_{\max} \simeq \lambda^{\star}_{\min} \approx 1$  for notational simplicity.

1. Linear convergence. In the absence of noise, the proposed algorithm converges linearly; namely, it provably attains  $\varepsilon$  accuracy within  $O(\log(1/\varepsilon))$  iterations. Given the inexpensiveness of each gradient iteration, this algorithm can be viewed as a linear-time algorithm, which can almost be implemented as long as we can

read the data. In the noisy setting, the algorithm reaches an appealing statistical accuracy within a logarithmic number of iterations.

- 2. Near-optimal sample complexity. The fast convergence is guaranteed as soon as the sample size exceeds the order of  $d^{3/2}$  poly log d. This matches the minimal sample complexity—modulo some logarithmic factor—known so far for any polynomial-time algorithm.
- 3. Near-optimal statistical accuracy. The proposed algorithm converges geometrically fast to a point with Euclidean error  $O(\sigma\sqrt{(d\log d)/p})$ . This matches the lower bound established in (Xia et al. 2017, theorem 5) up to some logarithmic factor, thus justifying the statistical optimality of the proposed nonconvex algorithm.
- 4. Entrywise estimation accuracy. In addition to the Euclidean statistical guarantees, we have also established an entrywise error bound, which, to the best of our knowledge, has not been established in any of the prior work. When t is sufficiently large, the iterates reach an entrywise error bound  $O(\sigma\sqrt{(\log d)/p})$ . This entrywise error bound is about an order of  $\sqrt{d}$  times smaller than the above  $\ell_2$  error bound, thereby implying that the estimation errors are evenly spread out across all entries.
- 5. Noise size. The aforementioned theory operates in the regime where  $\sigma \lesssim \sqrt{\frac{p}{d^{3/2}}}$  (modulo some log factor). Given that we have  $||T^{\star}||_{\infty} \asymp d^{-3/2}$  in this case, our noise size constraint can be equivalently written as (up to some log factor)

$$\frac{\sigma}{\|T^*\|_{\infty}} \lesssim \sqrt{pd^{3/2}}.\tag{15}$$

Since the sampling rate needs to satisfy  $p \gg d^{-3/2}$ , this condition essentially allows the typical size of each noise component to be considerably larger than the size of the corresponding entry of the truth, which covers a broad range of practical scenarios.

- 6. Implicit regularization. One appealing feature of our finding is the simplicity of the iterative refinement stage of the algorithm. All of the aforementioned statistical and computational benefits hold for vanilla gradient descent (when properly initialized). This should be contrasted with prior work (e.g., Xia and Yuan 2017) that relies on extra regularization terms to stabilize the optimization landscape. In principle, vanilla gradient descent implicitly constrains itself within a region of well-conditioned landscape, thus enabling fast convergence without explicit regularization.
- 7. No need of sample splitting. The theory developed herein does not require fresh samples in each iteration. We note that sample splitting has been frequently adopted in other context primarily to simplify mathematical analysis. Nevertheless, it typically does not exploit the data in an efficient manner (i.e., each data sample is used only once), thus resulting in the need of a much larger sample size in practice.

	Algorithm	Sample complexity	Computational complexity	$l_2$ error (noisy)	$l_{\infty}$ error (noisy)	Recovery type (noiseless)
Our theory	Spectral method + (vanilla) GD	$d^{1.5}$	pd <sup>3</sup>	$\sigma \sqrt{\frac{d}{p}}$	$\sigma\sqrt{\frac{1}{p}}$	Exact
Xia et al. (2017)	Spectral initialization + tensor power method	d <sup>1.5</sup>	$pd^3$	$(\ T^\star\ _{\infty} + \sigma)\sqrt{\frac{d}{p}}$	N/A	Approximate
Xia and Yuan (2017)	Spectral method + GD on manifold	$d^{1.5}$	poly(d)	N/A	N/A	Exact
Montanari and Sun (2018)	Spectral method	$d^{1.5}$	$d^3$	N/A	N/A	Approximate
Barak and Moitra (2016)	Sum-of-squares	$d^{1.5}$	$d^{15}$	$\frac{\ T^*\ _{\rm F}}{\sqrt{pd^{1.5}}} + \sigma d^{1.5}$	N/A	Approximate
Potechin and Steurer (2017)	Sum-of-squares	$d^{1.5}$	$d^{10}$	N/A	N/A	Exact
Yuan and Zhang (2016)	Tensor nuclear norm	d	NP-hard	N/A	N/A	Exact
Yuan and Zhang (2017)	Minimization					

**Table 1.** Comparison with Prior Theory for Existing Methods When  $r, \mu, \kappa \times 1$  (Neglecting Logarithmic Factors)

*Note*. N/A, not applicable.

We shall take a moment to discuss the merits of our approach in comparison with prior work. One of the best-known polynomial-time algorithms is the degree-6 level of the sum-of-squares (SoS) hierarchy, which seems to match the computationally feasible limit in terms of the sample complexity (Barak and Moitra 2016). However, this approach has a well-documented limitation in that it involves solving a semidefinite program of dimensions  $d^3 \times d^3$ , which requires enormous storage and computation power. The work of Montanari and Sun (2018) alleviates this computational burden by resorting to a clever unfolding-based spectral algorithm; it is a nearly linear-time procedure that enables near-minimal sample complexity (among polynomial-time algorithms), although it does not achieve exact recovery even in the absence of noise. The twostage algorithm developed by Xia et al. (2017)—which is based on spectral initialization followed by tensor power methods—shares similar advantages and drawbacks as the method developed by Montanari and Sun (2018). Further, the recent work of Xia and Yuan (2017) proposes a polynomial-time nonconvex algorithm based on gradient descent over a Grassmann manifold (with a properly regularized objective function), which is an extension of the nonconvex matrix completion algorithm proposed by (Keshavan et al. 2010a,b) to tensor data. The theory provided by Xia and Yuan (2017), however, does not provide explicit computational complexities. The recent work of Shah and Yu (2019) attempts tensor estimation via an interesting algorithm adapted from collaborative filtering and investigates both  $\ell_2$  and  $\ell_{\infty}$  estimation accuracy. This approach, however, does not guarantee exact recovery in the

absence of noise. We summarize and compare several prior results in Table 1 (omitting logarithmic factors).

Thus far, we have concentrated on the low-rank, well-conditioned, and incoherent case. Our main theory can be extended to cover a broader class of scenarios.

**Theorem 2.** Fix an arbitrary small constant  $\delta > 0$ . Suppose that  $\kappa \times 1$ ,

$$p \ge c_0 \frac{\mu^4 r^4 \log^4 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^*} \le c_1 \frac{\sqrt{p}}{\mu r^{3/2} d^{3/4} \log^2 d},$$

$$r \le c_2 \left(\frac{d}{\mu^6 \log^6 d}\right)^{1/6},$$

$$L = c_3 r^{2\kappa^2} \log^{3/2} r \quad \text{and}$$

$$\epsilon_{\text{th}} = c_4 \left(\frac{\mu r \log d}{d\sqrt{p}} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{r d \log^2 d}{p}} + \sqrt{\frac{\mu r \log d}{d}}\right)$$

for some sufficiently large constants  $c_0, c_3 > 0$  and some sufficiently small constants  $c_1, c_2, c_4 > 0$ . The learning rate  $\eta_t \equiv \eta$  is taken to be a constant obeying  $0 < \eta \leq \lambda_{\min}^{\star 4/3}/(32\lambda_{\max}^{\star 8/3})$ . Then with probability at least  $1 - \delta$ ,

$$\operatorname{dist}_{F}(\boldsymbol{U}^{t}, \boldsymbol{U}^{\star}) \leq \left(C_{1} \rho^{t} + C_{2} \frac{\sigma}{\lambda_{\min}^{\star}} \sqrt{\frac{d \log d}{p}}\right) \|\boldsymbol{U}^{\star}\|_{F}, \quad (16a)$$

$$\operatorname{dist}_{\infty}(\boldsymbol{U}^{t}, \boldsymbol{U}^{\star}) \leq \operatorname{dist}_{2,\infty}(\boldsymbol{U}^{t}, \boldsymbol{U}^{\star})$$

$$\leq \left(C_{3}\rho^{t} + C_{4}\frac{\sigma}{\lambda_{\min}^{\star}} \sqrt{\frac{d\log d}{p}}\right) \|\boldsymbol{U}^{\star}\|_{2,\infty} \tag{16b}$$

hold simultaneously for all  $0 \le t \le t_0 = d^5$ . Here,  $0 < C_1$ ,  $C_3$ ,  $\rho < 1$  and  $C_2$ ,  $C_4 > 0$  are some absolute constants.

**Corollary 2.** Fix an arbitrarily small constant  $\delta > 0$ . Instate the assumptions of Theorem 2. Then with probability at least  $1 - \delta$ ,

$$\| \mathbf{T}^t - \mathbf{T}^\star \|_{\mathrm{F}} \lesssim \left( C_1 \rho^t + C_2 \frac{\sigma}{\lambda_{\min}^\star} \sqrt{\frac{d \log d}{p}} \right) \| \mathbf{T}^\star \|_{\mathrm{F}}, \quad (17a)$$

$$\|T^{t} - T^{\star}\|_{\infty} \lesssim \left(C_{3}\rho^{t} + C_{4}\frac{\sigma}{\lambda_{\min}^{\star}}\sqrt{\frac{\mu^{3}rd\log d}{p}}\right)\|T^{\star}\|_{\infty} \quad (17b)$$

hold simultaneously for all  $0 \le t \le t_0 = d^5$ . Here,  $0 < C_1$ ,  $C_3$ ,  $\rho < 1$  and  $C_2$ ,  $C_4 > 0$  are some absolute constants.

**Remark 4.** Clearly, Theorem 2 and Corollary 2 subsume Theorem 1 and Corollary 1, respectively, as special cases.

**Remark 5.** Our theorems require the rank r to not exceed  $o(d^{1/6})$ , which, we believe, is an artifact of the current nonconvex analysis (particularly for the initialization stage). For instance, our local convergence analysis is built upon strong convexity and smoothness, which holds only within a sufficiently small neighborhood surrounding the truth; given that the diameter of this neighborhood is no more than o(1/r), our analysis requires an initial guess with higher accuracy than expected, thus leading to our rank constraint. It might be possible to improve the rank dependency via more refined analysis, and we leave it to future investigation.

In a nutshell, this theorem reveals intriguing theoretical support (including both  $\ell_{\rm F}$  and  $\ell_{2,\infty}$  bounds) for more general settings. Assuming that the condition number  $\kappa \approx 1$ , the nonconvex algorithm that we propose is guaranteed to succeed in polynomial time. Note, however, that our theoretical dependency (including both sample and computational complexities) on the rank r and the incoherence parameter  $\mu$  are likely loose and suboptimal. In addition, if  $\kappa$  is allowed to grow with d, then the current theory requires a large number of restart attempts during the initialization stage, resulting in a very high computational burden. Improving these aspects, however, calls for a much more refined analysis framework, which we leave for future investigation.

### 2.3. Numerical Experiments

We carry out a series of numerical experiments to corroborate our theoretical findings. Before proceeding, recall that Theorem 2 only guarantees successful recovery with probability  $1 - \delta$  for some small constant  $\delta$ ; this means that we shall not anticipate a very high success rate (e.g.,  $1 - O(d^{-5})$ ), as in the matrix recovery case. As we shall make clear shortly, this happens mainly because the initialization stage works only with probability  $1 - \delta$ , where the uncertainty largely

depends on the random vectors  $\{g^{\tau}\}_{1 \le \tau \le L}$ . With this observation in mind, we recommend the following modification to improve the empirical success rate:

• Run Algorithm 2 independently for  $t_{\text{init}} = 5$  times to obtain multiple initial estimates (denoted by  $U_{[1]}^0, \dots, U_{[t_{\text{init}}]}^0$ ); select the one achieving the smallest empirical loss, namely,

$$U_{\text{best}}^{0} = \underset{U \in \{U_{[i]}^{0}\}_{1 \le i \le h_{\text{nit}}}}{\text{arg min }} f(U).$$
 (18)

• Run Algorithm 1 with the initial point  $U^0$  set to be  $U^0_{\text{best}}$ .

The final estimates for the low-rank factor and the whole tensor are denoted, respectively, by

$$\hat{\mathbf{U}} = \mathbf{U}^{t_0} \quad \text{and} \quad \hat{T} = \sum_{i=1}^{r} \mathbf{u}_i^{t_0} \otimes \mathbf{u}_i^{t_0} \otimes \mathbf{u}_i^{t_0},$$
 (19)

where  $\boldsymbol{U}^{t_0} = [\boldsymbol{u}_1^{t_0}, \cdots, \boldsymbol{u}_r^{t_0}] \in \mathbb{R}^{d \times r}$  is the iterate returned by Algorithm 1, with  $t_0$  the total number of gradient iterations. In the sequel, we generate the true tensor  $T^\star = \sum_{1 \leq i \leq r} \boldsymbol{u}_i^{\star \otimes 3}$  randomly in such a way that  $\boldsymbol{u}_i^{\star \text{i.i.d.}} \sim \mathcal{N}(0, \boldsymbol{I}_d)$ . The learning rates are taken to be  $\eta_t \equiv 0.2$  unless otherwise noted.

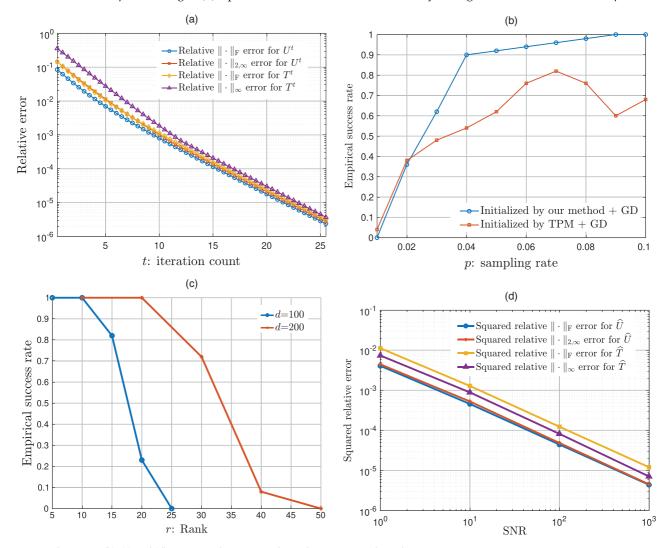
We start with numerical convergence rates of our algorithm in the absence of noise. Set d=100, r=4, p=0.1, L=16, and  $\epsilon_{th}=0.4$ . Figure 1(a) shows the numerical estimation errors versus iteration count t in a typical Monte Carlo trial. Here, four kinds of estimation errors are reported: (1) the relative Frobenius norm error  $\frac{\text{dist}_{E}(U^{t}, U^{*})}{\|U^{*}\|_{E}}$ ; (2) the relative  $\|\cdot\|_{2,\infty}$  error  $\frac{\text{dist}_{2,\infty}(U^{t}, U^{*})}{\|U^{*}\|_{2,\infty}}$ ; (3) the relative Frobenius norm error  $\frac{\|T^{t}-T^{*}\|_{E}}{\|T^{*}\|_{F}}$ ; and (4)

the relative Frobenius norm error  $\frac{u}{\|T^t\|_F}$ ; and (4) the relative  $\ell_{\infty}$  error  $\frac{\|T^t-T^t\|_{\infty}}{\|T^t\|_{\infty}}$ . Here,  $T^t = \sum_{i=1}^r u_i^t \otimes u_i^t \otimes u_i^t \otimes u_i^t$  with  $U^t = [u_1^t, \cdots, u_r^t]$ . For all these metrics, the numerical estimation errors decay geometrically fast.

Next, we study the phase transition (in terms of the success rates for exact recovery) in the noise-free settings. Set d=100, r=4, L=16, and  $\epsilon_{th}=0.4$ . For the sake of comparisons, we also report the numerical performance of the tensor power method (TPM) followed by gradient descent. When running the tensor power method, we set both the number of iterations and the restart number to be 16. Each trial is claimed to succeed if the relative  $\ell_2$  error obeys  $\frac{\text{dist}_F(\hat{\boldsymbol{U}}, \boldsymbol{U}^*)}{\|\boldsymbol{U}^*\|_F} \leq 0.01$ . Figure 1(b) plots the empirical success rates over 100 independent Monte Carlo trials. As can be seen, our initialization algorithm outperforms the tensor power method.

The third series of experiments is concerned with the dependence of the success rate on the rank r. Let us set  $p = rd^{-3/2}\log^2 d$ ,  $L = r^2$ , and  $\epsilon_{th} = 0.4$ , and the success recovery criterion is the same as mentioned earlier. Figure 1(c) depicts the empirical success rates (over 100 independent Monte Carlo trials) as the rank r varies. As can be seen from the plots, the proposed

**Figure 1.** (Color online) (a) Relative Errors of the Estimates  $U^t$  and  $T^t$  vs. Iteration Count t for Noiseless Tensor Completion, Where d = 100, r = 4, and p = 0.1; (b) Empirical Success Rate vs. Sampling Rate, Where d = 100 and r = 4; (c) Empirical Success Rate vs. Rank, Where  $p = rd^{-3/2}\log^2 d$ ; (d) Squared Relative Errors vs. SNR for Noisy Settings, Where d = 100, r = 4, and p = 0.1



Note. Each point in (b), (c) and (d) is averaged over 100 independent Monte Carlo trials.

algorithm is able to achieve exact reconstruction as long as the rank r is sufficiently small compared with d. The plausible range of r, however, seems to be larger than our theoretic requirement  $r = o(d^{1/6})$ . This, once again, suggests the need of future investigation to pin down the best possible dependency on r.

Finally, we consider the numerical estimation accuracy of our algorithm. Take  $t_0=100$ , d=100, r=4, p=0.1, L=16, and  $\epsilon_{\rm th}=0.4$ . Define the signal-to-noise ratio (SNR) to be  ${\rm SNR}=\frac{\|T^*\|_{\rm F}^2/d^3}{\sigma^2}$ . We report in Figure 1(d) three types of squared relative errors (i.e.,  $\frac{{\rm dist}_{\rm F}^2(\hat{\boldsymbol{U}},\boldsymbol{U}^*)}{\|\boldsymbol{U}^*\|_{\rm F}^2}$ ,  $\frac{{\rm dist}_{2,\infty}^2(\hat{\boldsymbol{U}},\boldsymbol{U}^*)}{\|\boldsymbol{U}^*\|_{2,\infty}^2}$ , and  $\frac{\|\hat{\boldsymbol{T}}-T^*\|_{\infty}^2}{\|T^*\|_{\infty}^2}$ ) versus SNR. Figure 1(d) illustrates that all three types of relative squared errors scale inversely proportional to the SNR (since the slope in the figure is roughly -1), which is

consistent with our statistical guarantees.

#### 2.4. Notation

Before proceeding, we gather a few notations that will be used throughout this paper. First of all, for any matrix  $M \in \mathbb{R}^{d \times d}$ , we let  $\|M\|$  and  $\|M\|_F$  denote the operator norm (or the spectral norm) and the Frobenius norm of M, respectively, and we let  $M_{i,:}$  and  $M_{:,i}$  denote the i th row and i th column, respectively. In addition, we let  $\lambda_1(M) \geq \lambda_2(M) \geq \cdots \geq \lambda_d(M)$  denote the eigenvalues of M and  $\sigma_1(M) \geq \sigma_2(M) \geq \cdots \geq \sigma_d(M)$  denote the singular values of M.

note the singular values of M. For any tensor  $T \in \mathbb{R}^{d \times d \times d}$ , let  $T_{i,:,:} \in \mathbb{R}^{d \times d}$  denote the mode-1 i-slice with entries  $(T_{i,:,:})_{j,k} = T_{i,j,k}$ ;  $T_{:,i,:}$  and  $T_{:,:,:}$  are defined in a similar way. For any tensors  $T, R \in \mathbb{R}^{d \times d \times d}$ , the inner product is defined as  $\langle T, R \rangle := \sum_{i,j,k,l} T_{j,k,l} R_{j,k,l}$ . The Frobenius norm of T is defined as  $\|T\|_F := \sqrt{\langle T, T \rangle}$ . For any vectors  $u, v \in \mathbb{R}^d$ , we define the vector products of a tensor  $T \in \mathbb{R}^{d \times d \times d}$ —denoted by  $T \times_3 u \in \mathbb{R}^{d \times d}$  and  $T \times_1 u \times_2 v \in \mathbb{R}^d$ —such that

$$[T \times_3 u]_{ij} := \sum_{1 \le k \le d} T_{i,j,k} u_k, \quad 1 \le i, j \le d;$$
 (20a)

$$[T \times_1 u \times_2 v]_k := \sum_{1 \le i, j \le d} T_{i,j,k} u_i v_j, \quad 1 \le k \le d.$$
 (20b)

The products  $T \times_2 u \in \mathbb{R}^{d \times d}$ ,  $T \times_3 u \in \mathbb{R}^{d \times d}$ ,  $T \times_1 u \times_3 v \in \mathbb{R}^d$ , and  $T \times_2 u \times_3 v \in \mathbb{R}^d$  are defined in a similar manner. For any  $U = [u_1, \cdots, u_r] \in \mathbb{R}^{d \times r}$  and  $V = [v_1, \cdots, v_r] \in \mathbb{R}^{d \times r}$ , we further define

$$T \times_{1}^{\text{seq}} U \times_{2}^{\text{seq}} V := [T \times_{1} u_{i} \times_{2} v_{i}]_{1 \le i \le r} \in \mathbb{R}^{d \times r}.$$
 (21)

In addition, the operator norm of *T* is defined as

$$||T|| := \sup_{u,v,w \in \mathbb{S}^{d-1}} \langle T, u \otimes v \otimes w \rangle, \tag{22}$$

where  $\mathbb{S}^{d-1} := \{ u \in \mathbb{R}^d | ||u||_2 = 1 \}$  indicates the unit sphere in  $\mathbb{R}^d$ .

Further,  $f(n) \leq g(n)$  or f(n) = O(g(n)) means that  $|f(n)/g(n)| \leq C_1$  for some constant  $C_1 > 0$ ;  $f(n) \geq g(n)$  means that  $|f(n)/g(n)| \geq C_2$  for some constant  $C_2 > 0$ ;  $f(n) \approx g(n)$  means that  $C_1 \leq |f(n)/g(n)| \leq C_2$  for some constants  $C_1, C_2 > 0$ ; and f(n) = o(g(n)) means that  $\lim_{n \to \infty} f(n)/g(n) = 0$ . In addition,  $f(n) \ll g(n)$  means that  $f(n) \leq c_1g(n)$  for some sufficiently small constant  $c_1 > 0$ , and  $f(n) \gg g(n)$  means that  $f(n) \geq c_2g(n)$  for some sufficiently large constant  $c_2 > 0$ .

#### 3. Initialization

This section presents formal details of the proposed two-step initialization, accompanied by some intuition. We defer the discussion about alternative approaches to the supplementray materials. Recall that the proposed initialization procedure consists of two steps, which we discuss separately.

# 3.1. Step 1: Subspace Estimation via a Spectral Method

The spectral algorithm is often applied in conjunction with simple "unfolding" (or "matricization") to estimate the *subspace* spanned by the r factors  $\{u_i^\star\}_{1\leq i\leq r}$ . This strategy is partly motivated by prior approaches developed for covariance estimation with missing data (Lounici 2014, Montanari and Sun 2018, Cai et al. 2021). We next provide a brief introduction.

Let

$$A = \text{unfold}^{1\times 2} \left(\frac{1}{p}T\right) \in \mathbb{R}^{d\times d^2}, \text{ or more concisely}$$

$$A = \text{unfold}\left(\frac{1}{p}T\right) \in \mathbb{R}^{d\times d^2}$$
(23)

be the mode-1 matricization of  $p^{-1}T$  (i.e.,  $\frac{1}{p}T_{i,j,k} = A_{i,(j-1)d+k}$  for any  $1 \le i,j,k \le d$ ) (Kolda and Bader 2009).

The rationale of this step is that, under our model, the unfolded matrix *A* obeys

$$\mathbb{E}[A] = \mathsf{unfold}(T^*) = \sum_{i=1}^r u_i^* (u_i^* \otimes u_i^*)^\top =: A^*, \tag{24}$$

whose column space is precisely the span of  $\{u^*\}_{1 \le i \le r}$ . This motivates one to estimate the r-dimensional column space of  $\mathbb{E}[A]$  from A. Toward this, a natural strategy is to look at the principal subspace of  $AA^\top$ . However, the diagonal entries of  $AA^\top$  bear too much influence on the principal directions and need to be properly down-weighed. The current paper chooses to work with the principal subspace of the following matrix that zeros out all diagonal components:

$$B := \mathcal{P}_{\mathsf{off-diag}}(AA^{\mathsf{T}}),\tag{25}$$

where  $\mathcal{P}_{\mathsf{off-diag}}(\mathbf{Z})$  extracts out the off-diagonal entries of a squared matrix  $\mathbf{Z}$ . If we let  $\mathbf{U} \in \mathbb{R}^{d \times r}$  be an orthonormal matrix whose columns are the top-r eigenvectors of  $\mathbf{B}$ , then  $\mathbf{U}$  serves as our subspace estimate. See Algorithm 2 for a summary of the procedure.

# 3.2. Step 2: Retrieval of Low-Rank Tensor Factors from the Subspace Estimate

**3.2.1. Procedure.** As it turns out, it is possible to obtain rough (but reasonable) estimates of all individual low-rank tensor factors  $\{u_i^*\}_{1 \le i \le r}$ —up to global permutation—given a reliable subspace estimate U. This is in stark contrast to the low-rank matrix recovery case, where there exists some global rotational ambiguity that prevents us from disentangling the r factors of interest.

We begin by describing how to retrieve *one* tensor factor from the subspace estimate—a procedure summarized in Retrieve-one-tensor-factor(). Let us generate a random vector from the provided subspace *U* (which has orthonormal columns), that is,

$$\theta = \underbrace{UU^{\top}g}_{\text{projection of }g}$$
,  $g \sim \mathcal{N}(0, I_d)$ . (26)

The rescaled tensor data  $p^{-1}T$  is then transformed into a matrix via proper "projection" along this random direction  $\theta$ , namely,

$$M = \frac{1}{p} T \times_3 \theta \in \mathbb{R}^{d \times d}.$$
 (27)

Our estimate for a tensor factor is then given by  $\lambda^{1/3}\nu$ , where  $\nu$  is the leading singular vector of M obeying  $\langle T, \nu^{\otimes 3} \rangle \geq 0$ , and  $\lambda$  is taken as  $\lambda = \langle p^{-1}T, \nu^{\otimes 3} \rangle$ . Informally,  $\nu$  reflects the direction of the component  $u_i^*$  that exhibits the largest correlation with the random direction  $\theta$ , and  $\lambda$  forms an estimate of the corresponding size  $\|u_i^*\|_2$ .

A challenge remains, however, as there are oftentimes more than one tensor factor to estimate. To address this issue, we propose to rerun the aforementioned procedure multiple times, so as to ensure that we get to retrieve each tensor factor of interest at least once. We will then apply a careful pruning procedure (i.e., PRUNE()) to remove redundancy.

**3.2.2. Intuition.** To develop some intuition about the aforementioned procedure, consider the "heuristic" case where  $\theta = U^*(U^{*\top}U^*)^{-1}U^{*\top}g$ , namely, the idealistic scenario where the subspace estimate U is accurate. Averaging out the randomness in the sampling pattern and the noise, we see that the expected projected matrix (27) takes the following form:

$$\mathbb{E}[M | \theta] = T^* \times_3 \theta = \sum_{i=1}^r \langle \theta, u_i^* \rangle u_i^* u_i^{*\top}.$$

As a result, in the incoherent case where  $\{u_j^\star\}$  are nearly orthogonal to each other, the leading singular vector of  $\mathbb{E}[M|\theta]$ —and hence that of M (i.e., w)—is expected to be reasonably close to the factor  $u_i^\star$  that enjoys the largest projected coefficient. In other words, we expect

$$v \approx \frac{1}{\|\boldsymbol{u}_{i}^{\star}\|_{2}} \boldsymbol{u}_{i}^{\star}, \text{ where } i = \arg\max_{1 \leq j \leq r} |\langle \boldsymbol{\theta}, \boldsymbol{u}_{j}^{\star} \rangle|.$$
 (28)

In the mean time, armed with (28) and the incoherence assumption (such that  $u_i^*$  and  $u_j^*$  are nearly orthogonal for  $i \neq j$ ), one might have

$$\lambda = \langle T^{\star}, \boldsymbol{\nu}^{\otimes 3} \rangle \approx \frac{1}{\parallel \boldsymbol{u}_{i}^{\star} \parallel_{2}^{3}} \langle T^{\star}, \boldsymbol{u}_{i}^{\star \otimes 3} \rangle$$

$$\approx \frac{1}{\parallel \boldsymbol{u}_{i}^{\star} \parallel_{2}^{3}} \langle \boldsymbol{u}_{i}^{\star \otimes 3}, \boldsymbol{u}_{i}^{\star \otimes 3} \rangle = \parallel \boldsymbol{u}_{i}^{\star} \parallel_{2}^{3},$$
(29)

thus explaining our choice of  $\lambda$  in the proposed procedure. These arguments hint at the ability of our procedure in retrieving one tensor factor in each round.

This intuitive argument, however, does not explain why we need to first project a random vector g onto the (approximate) column space of  $U^*$ . Although we will not go into detailed calculations here, we remark in passing a crucial high variability issue: without proper projection, the perturbation incurred by both the missing data and the noise might far exceed the strength of the true signal. As a result, it is advised to first project the data onto the desired subspace, in the hope of amplifying the signal-to-noise ratio.

### 4. Related Work

One of the most natural ideas for solving tensor completion is to first unfold the tensor data into matrices, followed by proper convex relaxation commonly adopted for low-rank matrix completion. Given that there is more than one way to matricize a tensor, several prior works have explored the design of matrix norms that can exploit the tensor structure more effectively (Tomioka et al. 2010, Gandy et al. 2011, Liu et al. 2013, Romera-Paredes and Pontil 2013, Mu et al. 2014, Lu et al. 2016). Such algorithms have been robustified

to enable reliable recovery against sparse outliers as well (Goldfarb and Qin 2014). For the most part, however, such unfolding-based convex relaxation necessarily incurs loss of structural information, which is particularly severe when handling odd-order tensors. The sample complexity developed for this paradigm is often suboptimal vis-à-vis the computational limits (i.e., minimal sample complexity achievable by polynomial-time algorithms).

Motivated by the aforementioned suboptimality issue, Yuan and Zhang (2016, 2017) proposed to minimize instead the tensor nuclear norm subject to data constraints, which provably allows for reduced sample complexity. The issue, however, is that computing the tensor nuclear norm itself is already computationally intractable, thus limiting its applicability to even moderate-dimensional problems. Similar findings have also been discovered for tensor atomic norm minimization (Driggs et al. 2019). When restricted to polynomial-time algorithms, the best statistical guarantees are often attained via convex relaxation tailored to the sum-of-squares hierarchy (Barak and Moitra 2016); the resulting computational cost, however, remains prohibitively high for practical large-scale problems. Another matrix nuclear norm minimization algorithm has been proposed based on promoting certain structures on certain factor matrices (Liu et al. 2014). Developing statistical guarantees is, however, not the focal point of this work.

Moving beyond convex relaxation, a number of prior papers have developed nonconvex algorithms for tensor completion, examples including iterative hard thresholding (Rauhut et al. 2017), alternating minimization (Jain and Oh 2014, Xu et al. 2015, Wang et al. 2016), tensor SVD (Zhang and Aeron 2017), optimization on manifold (Kasai and Mishra 2016, Steinlechner 2016, Xia and Yuan 2017), proximal average algorithm with nonconvex regularizer (Yao 2018), and block coordinate decent (Xu and Yin 2013, Ji et al. 2016). When it comes to the model considered herein, these algorithms either lack optimal statistical guarantees or come with a computational cost that is significantly higher than a linear-time algorithm.

The algorithm and theory that we develop are largely inspired by the recent advances of nonconvex optimization algorithms for low-rank matrix recovery problems (Keshavan et al. 2010a,b; Candès et al. 2015; Chen and Wainwright 2015; Sun and Luo 2016; Yi et al. 2016; Chen and Candès 2017). The main theoretical tool—the leave-one-out analysis—is a powerful technique that has proved successful in various other statistical problems (El Karoui 2015, Abbe et al. 2017, Ding and Chen 2018, Zhong and Boumal 2018, Chen et al. 2019c, Chen et al. 2019d, Chen et al. 2019e, Li et al. 2019, Pananjady and Wainwright 2019, Ma et al. 2020, Chen et al. 2021). There are several major differences

between the analysis of nonconvex tensor completion and that of nonconvex matrix recovery. For instance, our initialization scheme is substantially more complicated than the matrix recovery counterpart, thus requiring much more sophisticated analysis; in addition, the local convergence stage of tensor completion does not suffer from rotational ambiguity (which often appears in nonconvex matrix completion), and hence we only need to handle permutational ambiguity.

In addition, the current paper focuses on nonadaptive uniform random sampling. If there is freedom in designing the sampling mechanism, then one can often expect improved performance (see, e.g., Krishnamurthy and Singh 2013, Zhang 2019). Fundamental criteria that enable perfect low-CP-rank tensor completion have been studied by Ashraphijuo and Wang (2017).

Tensor completion is simply a special example of the tensor recovery literature. There is a large body of results tackling various other tensor recovery and estimation problems, including, but not limited to, tensor decomposition (Kolda 2001; Kolda and Bader 2009; Kim et al. 2013; Anandkumar et al. 2014a, b; Ge et al. 2015; Tang and Shah 2015; Hopkins et al. 2016; Ge and Ma 2017; Sidiropoulos et al. 2017; Sun et al. 2017; Zoubir et al. 2018), tensor SVD and factorization (Kilmer et al. 2013, Zhang and Aeron 2017, Zhang and Xia 2018), and tensor regression and sketching (Rauhut et al. 2017, Chen et al. 2019b, Hao et al. 2019, Hao et al. 2020). The algorithmic ideas explored in this paper might have implications for these tensor-related problems as well.

Finally, we remark that, compared with the conference version (Cai et al. 2019), the current paper (1) extends the results presented therein to a more general case, where both the rank r and the incoherence parameter  $\mu$  are allowed to grow with d; (2) discusses how to handle asymmetric tensors; and (3) explains in detail the inadequacy of other initialization schemes (including both random initialization and tensor power methods). More numerical experiments have also been carried out and reported.

# 5. Analysis

In this section, we outline the proof of Theorem 2 and defer the detailed proof to the e-companion. The analysis is divided into three parts:

- In Section 5.1, we show that, given an initial estimate sufficiently close to the ground truth, vanilla gradient descent converges linearly. These are formalized in Lemmas 3 and 6.
- Sections 5.2–5.3 provide statistical guarantees for the two steps of the initialization procedure (see Theorem 3).
- Under the assumptions of Theorem 2, one can see that the initialization satisfies the requirement of linear convergence of vanilla gradient descent. Therefore,

Theorem 2 immediately follows from the results in Sections 5.1–5.3.

## 5.1. Analysis for Local Convergence of GD

In this section, we demonstrate that if the initialization is reasonably good, then vanilla gradient descent converges linearly to a solution with the desired statistical accuracy. We postpone the analysis for initialization to Sections 5.2–5.3 for convenience of presentation.

## 5.1.1. Preliminaries: Gradient and Hessian Calculation.

First of all, using our notation  $x^{seq}$  defined in (21), we can write

$$\nabla f(\mathbf{U}) = \frac{1}{p} \mathcal{P}_{\Omega} \left( \sum_{1 \le i \le r} \mathbf{u}_i^{\otimes 3} - \mathbf{T}^* - \mathbf{E} \right) \times_1^{\mathsf{seq}} \mathbf{U} \times_2^{\mathsf{seq}} \mathbf{U}. \tag{30}$$

Next, we find it convenient to define an auxiliary loss function  $f_{\text{clean}}(\boldsymbol{U}): \mathbb{R}^{d \times r} \to \mathbb{R}_+$  that corresponds to the noiseless case:

$$f_{\text{clean}}(\boldsymbol{U}) = \frac{1}{6p} \left\| \mathcal{P}_{\Omega} \left( \sum_{1 \le i \le r} \boldsymbol{u}_{i}^{\otimes 3} - \boldsymbol{T}^{\star} \right) \right\|_{F}^{2}.$$
 (31)

The gradient of  $f_{\text{clean}}$  with respect to (w.r.t.)  $u_s$  ( $1 \le s \le r$ ) is thus given by

$$\nabla_{\boldsymbol{u}_{s}} f_{\text{clean}}(\boldsymbol{U}) = \frac{1}{p} \mathcal{P}_{\Omega} \left( \sum_{1 \leq i \leq r} \boldsymbol{u}_{i}^{\otimes 3} - \boldsymbol{T}^{\star} \right) \times_{1} \boldsymbol{u}_{s} \times_{2} \boldsymbol{u}_{s}, \quad 1 \leq s \leq r,$$
(32)

and hence one can write

$$\nabla f_{\mathsf{clean}}(\boldsymbol{U}) = \frac{1}{p} \mathcal{P}_{\Omega} \left( \sum_{1 \le i \le r} \boldsymbol{u}_{i}^{\otimes 3} - \boldsymbol{T}^{\star} \right) \times_{1}^{\mathsf{seq}} \boldsymbol{U} \times_{2}^{\mathsf{seq}} \boldsymbol{U}. \tag{33}$$

This clearly satisfies

$$\nabla f(\mathbf{U}) = \nabla f_{\text{clean}}(\mathbf{U}) - \frac{1}{p} \mathcal{P}_{\Omega}(\mathbf{E}) \times_{1}^{\text{seq}} \mathbf{U} \times_{2}^{\text{seq}} \mathbf{U}. \tag{34}$$

Moreover, direct algebraic manipulations give that, for any matrix  $V = [v_1, ..., v_r] \in \mathbb{R}^{d \times r}$ ,

$$\operatorname{vec}(V)^{\top} \nabla^{2} f_{\operatorname{clean}}(U) \operatorname{vec}(V)$$

$$= \frac{1}{3p} \left\| \mathcal{P}_{\Omega} \left( \sum_{1 \leq s \leq r} u_{s} \otimes u_{s} \otimes v_{s} + u_{s} \otimes v_{s} \otimes u_{s} \right) + v_{s} \otimes u_{s} \otimes u_{s} \right\|_{F}^{2} + \frac{2}{p} \left\langle \mathcal{P}_{\Omega} \left( \sum_{s \in [r]} u_{s}^{\otimes 3} - T^{\star} \right) \right\rangle,$$

$$\times \sum_{s \in [r]} v_{s} \otimes v_{s} \otimes u_{s} \right\rangle,$$
(35)

where vec(V) denotes the vectorization of V.

**5.1.2. Local Strong Convexity and Smoothness.** At the heart of our analysis is a crucial geometric property of the objective function; that is, the noiseless loss function  $f_{\text{clean}}$  behaves like a locally strongly convex and smooth function. This fact, which is formally

stated in the following lemma, is the key enabler of fast local convergence of vanilla GD.

Lemma 1 (Local Strong Convexity and Smoothness).

$$p \ge c_0 \max\left\{\frac{\log^3 d}{d^{3/2}}, \frac{\mu^2 r^2 \log d}{d^2}\right\} \qquad r \le c_1 \sqrt{\frac{d}{\mu}}$$
 (36)

for some sufficiently large (respectively, small) constant  $c_0 > 0$  (respectively,  $c_1 > 0$ ). Then with probability greater than  $1 - O(d^{-10})$ ,

$$\frac{1}{2} \lambda_{\min}^{*4/3} \| V \|_{F}^{2} \le \text{vec}(V)^{\top} \nabla^{2} f_{\text{clean}}(U) \text{vec}(V) \le 4 \lambda_{\max}^{*4/3} \| V \|_{F}^{2}$$
(37)

holds simultaneously for all  $V \in \mathbb{R}^{d \times r}$  and all  $U \in \mathbb{R}^{d \times r}$  obeying

$$\|\boldsymbol{U} - \boldsymbol{U}^{\star}\|_{F} \le \delta \|\boldsymbol{U}^{\star}\|_{F} \quad and \quad \|\boldsymbol{U} - \boldsymbol{U}^{\star}\|_{2,\infty} \le \delta \|\boldsymbol{U}^{\star}\|_{2,\infty}.$$
(38)

Here,  $\delta \le c_2/(\mu^{3/2}r)$  for some sufficiently small constant  $c_2 > 0$ .

In order to invoke Lemma 1, one needs to make sure that the decision matrix  $\boldsymbol{U}$  of interest (e.g.,  $\boldsymbol{U}^t$  in the GD sequence) satisfies the condition (38). This, however, is a fairly stringent condition, as it requires  $\boldsymbol{U}$  to be close to the truth in every single row.

### 5.1.3. Leave-One-Out Gradient Descent Sequences.

Motivated by the analytical framework developed for low-rank matrix recovery (Ma et al. 2017, Chen et al. 2019a), we introduce the following leave-one-out sequences, which play a crucial role in guaranteeing that the entire trajectory  $\{\boldsymbol{U}^t\}_{t\geq 0}$  satisfies the condition (38), as required in Lemma 1.

Specifically, we define for each  $1 \le m \le d$  the following auxiliary loss function:

$$f^{(m)}(\mathbf{U}) \triangleq \frac{1}{6p} \left\| \mathcal{P}_{\Omega_{-m}} \left( \sum_{1 \leq s \leq r} \mathbf{u}_s^{\otimes 3} - \mathbf{T}^* - \mathbf{E} \right) \right\|_{\mathrm{F}}^2$$

$$+ \frac{1}{6} \left\| \mathcal{P}_m \left( \sum_{1 \leq s \leq r} \mathbf{u}_s^{\otimes 3} - \mathbf{T}^* \right) \right\|_{\mathrm{F}'}^2$$
(39)

where

- $\mathcal{P}_{\Omega_m}$  is the projection onto the subspace of tensors supported on  $\{(i,j,k) \in \Omega : i = m \text{ or } j = m \text{ or } k = m\};$
- $\mathcal{P}_{\Omega_{-m}}$  is the projection onto the subspace of tensors supported on  $\{(i,j,k) \in \Omega : i \neq m \text{ and } j \neq m \text{ and } k \neq m\}$ ;
- $\mathcal{P}_m$  is the projection onto the subspace of tensors supported on  $\{(i,j,k) \in [d]^3 : i = m \text{ or } j = m \text{ or } k = m\}$ .

In words, this function is obtained by replacing all data at locations  $\{(i,j,k) \in [d]^3 : i = m \text{ or } j = m \text{ or } k = m\}$  by their expected values, thus removing all randomness associated with this location

subset. The gradient of  $f^{(m)}(U)$  w.r.t.  $u_s$   $(1 \le s \le r)$  can be computed as

$$\nabla_{\mathbf{u}_{s}} f^{(m)}(\mathbf{U}) = \frac{1}{p} \mathcal{P}_{\Omega_{-m}} \left( \sum_{1 \leq s \leq r} \mathbf{u}_{s}^{\otimes 3} - \mathbf{T}^{\star} - \mathbf{E} \right)$$

$$\times_{1} \mathbf{u}_{s} \times_{2} \mathbf{u}_{s} + \mathcal{P}_{m} \left( \sum_{1 \leq s \leq r} \mathbf{u}_{s}^{\otimes 3} - \mathbf{T}^{\star} \right)$$

$$\times_{1} \mathbf{u}_{s} \times_{2} \mathbf{u}_{s}.$$

$$(40)$$

We then denote by  $\{\boldsymbol{U}^{t,(m)}\}_{t\geq 0}$  the iterative sequence obtained by running gradient descent w.r.t. the leave-one-out loss  $f^{(m)}(\cdot)$  (see Algorithm 4). By construction, as long as  $\boldsymbol{U}^{0,(m)}$  is independent of the sampling locations and the noise associated with the locations  $\{(i,j,k)\in\Omega:i=m\text{ or }j=m\text{ or }k=m\}$  (which holds true as detailed momentarily), then the entire trajectory  $\{\boldsymbol{U}^{t,(m)}\}_{t\geq 0}$  becomes statistically independent of such randomness. This is a crucial property that allows us to decouple the complicated statistical dependency.

**Algorithm 4** (The *m* th Leave-One-Out Sequence)

1: Generate an initial estimate  $U^{0,(m)}$  via Algorithm 5.

2: **for** 
$$t = 0, 1, ..., t_0 - 1$$
 **do**  
3:  $\mathbf{U}^{t+1,(m)} = \mathbf{U}^{t,(m)} - \eta_t \nabla f^{(m)} (\mathbf{U}^{t,(m)})$ .

**5.1.4. Key Lemmas.** The proof for local linear convergence of GD is inductive in nature, which proceeds on the basis of the following set of inductive hypotheses. As we shall see in Corollary 3 in Section 5.3, this set of inductive hypotheses—modulo some global permutation—is valid with high probability when t=0. In order to simplify presentation, we remove the consideration of the global permutation factor throughout this section (i.e., we assume that the following holds for  $U^0\Pi^0$  with some permutation matrix  $\Pi^0 \in \mathbb{R}^{r \times r}$  obeying  $\Pi^0 = I$ . Our key inductive hypotheses for the gradient update stage are summarized as follows:

$$\|\boldsymbol{U}^{t} - \boldsymbol{U}^{\star}\|_{F} \leq \left(C_{1} \rho^{t} \mathcal{E}_{\text{local}} + C_{2} \frac{\sigma}{\lambda_{\min}^{\star}} \sqrt{\frac{d \log d}{p}}\right) \|\boldsymbol{U}^{\star}\|_{F}; \quad (41a)$$

$$\|\boldsymbol{U}^{t} - \boldsymbol{U}^{\star}\|_{2,\infty} \leq \left(C_{3} \rho^{t} \mathcal{E}_{\text{local}} + C_{4} \frac{\sigma}{\lambda_{\min}^{\star}} \sqrt{\frac{d \log d}{p}}\right) \qquad (41b)$$
$$\times \|\boldsymbol{U}^{\star}\|_{2,\infty};$$

$$\|\boldsymbol{U}^{t} - \boldsymbol{U}^{t,(m)}\|_{F} \leq \left(C_{5}\rho^{t}\mathcal{E}_{local} + C_{6}\frac{\sigma}{\lambda_{\min}^{\star}}\sqrt{\frac{d\log d}{p}}\right) \qquad (41c)$$

$$\times \|\boldsymbol{U}^{\star}\|_{2,\infty};$$

$$\left\| \left( \boldsymbol{U}^{t,(m)} - \boldsymbol{U}^{\star} \right)_{m,:} \right\|_{2} \le \left( C_{7} \rho^{t} \mathcal{E}_{\text{local}} + C_{8} \frac{\sigma}{\lambda_{\min}^{\star}} \sqrt{\frac{d \log d}{p}} \right)$$
(41d)
$$\times \left\| \boldsymbol{U}^{\star} \right\|_{2,\infty};$$

for some quantity  $\mathcal{E}_{local} > 0$  (depending possibly on  $\mu$  and r) and some constants  $C_1, \dots, C_8 > 0$ . There exist a few straightforward consequences of the hypotheses (41), which we record in the following lemma.

**Lemma 2.** Assume that the hypotheses (41) hold. Then we have

$$\|\boldsymbol{U}^{t,(m)} - \boldsymbol{U}^{\star}\|_{F} \leq \left(2C_{1}\rho^{t}\mathcal{E}_{local} + 2C_{2}\frac{\sigma}{\lambda_{\min}^{\star}}\sqrt{\frac{d\log d}{p}}\right) \times \|\boldsymbol{U}^{\star}\|_{F},$$
(42)

$$\|\boldsymbol{U}^{t,(m)} - \boldsymbol{U}^{\star}\|_{2,\infty} \leq \left( (C_3 + C_5) \rho^t \mathcal{E}_{\text{local}} + (C_4 + C_6) \right.$$

$$\times \frac{\sigma}{\lambda_{\min}^{\star}} \sqrt{\frac{d \log d}{p}} \|\boldsymbol{U}^{\star}\|_{2,\infty}. \tag{43}$$

Our proof for the hypotheses (41) is inductive in nature: we would like to show that if the hypotheses in (41) hold for the t th iteration, then they continue to be valid for the (t+1) th iteration. We shall justify each of the aforementioned hypotheses inductively through the following lemmas.

Lemma 3. Suppose that

$$p \ge c_0 \frac{\mu^3 r^2 \log^3 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^*} \le c_1 \frac{\sqrt{p}}{\mu^{3/2} r \sqrt{d \log d}}, \quad \text{and}$$

$$r \le c_2 \sqrt{\frac{d}{\mu}}$$

for some sufficiently large constant  $c_0 > 0$  and some sufficiently small constant  $c_1, c_2 > 0$ . Assume that the hypotheses (41) hold for the tth iteration and  $\mathcal{E}_{local} \leq c_3/(\mu^{3/2}r)$  for some sufficiently small constant  $c_3 > 0$ . Then with probability at least  $1 - O(d^{-10})$ ,

$$\| \mathbf{U}^{t+1} - \mathbf{U}^{\star} \|_{F}$$

$$\leq \left( C_{1} \rho^{t+1} \mathcal{E}_{\text{local}} + C_{2} \frac{\sigma}{\lambda_{\min}^{\star}} \sqrt{\frac{d \log d}{p}} \right) \times \| \mathbf{U}^{\star} \|_{F},$$

$$(44)$$

provided that  $0 < \eta \le \lambda_{\min}^{\star 4/3}/(32\lambda_{\max}^{\star 8/3})$ ,  $1 - (\lambda_{\min}^{\star 4/3}/5)\eta \le \rho < 1$ , and  $C_2$  is sufficiently large.

Lemma 4. Suppose that

$$p \ge c_0 \frac{\mu^3 r^2 \log^3 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^{\star}} \le c_1 \frac{\sqrt{p}}{\mu^{3/2} r \sqrt{d \log d}}, \quad r \le c_2 \sqrt{\frac{d}{\mu}}$$

for some sufficiently large constant  $c_0 > 0$  and some sufficiently small constant  $c_1, c_2 > 0$ . Assume that the hypotheses (41) hold for the tth iteration and  $\mathcal{E}_{local} \leq c_3/(\mu^{3/2}r)$  for some sufficiently small constant  $c_3 > 0$ . Then with probability at least  $1 - O(d^{-10})$ , one has

$$\|\boldsymbol{U}^{t+1,(m)} - \boldsymbol{U}^{t+1}\|_{F} \le \left(C_{5}\rho^{t+1}\mathcal{E}_{local} + C_{6}\frac{\sigma}{\lambda_{\min}^{\star}}\sqrt{\frac{d\log d}{p}}\right) \times \|\boldsymbol{U}^{\star}\|_{2,\infty}, \tag{45}$$

provided that  $0 < \eta \le \lambda_{\min}^{\star 4/3}/(32\lambda_{\max}^{\star 8/3})$ ,  $1 - (\lambda_{\min}^{\star 4/3}/5)\eta \le \rho < 1$ , and  $C_6$  is sufficiently large.

**Lemma 5.** Suppose that

$$p \geq c_0 \frac{\mu^3 r^2 \mathrm{log}^3 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^\star} \leq c_1 \frac{\sqrt{p}}{\mu^{3/2} r \sqrt{d \log d}}, \quad r \leq c_2 \sqrt{\frac{d}{\mu}}$$

for some sufficiently large constant  $c_0 > 0$  and some sufficiently small constant  $c_1, c_2 > 0$ . Assume that the hypotheses (41) hold for the t th iteration and  $\mathcal{E}_{local} \leq c_3/(\mu^{3/2}r)$  for some sufficiently small constant  $c_3 > 0$ . Then with probability at least  $1 - O(d^{-10})$ , one has

$$\left\| \left( \boldsymbol{U}^{t+1,(m)} - \boldsymbol{U}^{\star} \right)_{m,:} \right\|_{2}$$

$$\leq \left( C_{7} \rho^{t+1} \mathcal{E}_{\text{local}} + C_{8} \frac{\sigma}{\lambda_{\min}^{\star}} \sqrt{\frac{d \log d}{p}} \right) \|\boldsymbol{U}^{\star}\|_{2,\infty}, \quad (46)$$

provided that  $0 < \eta \le \lambda_{\min}^{\star 4/3}/(32\lambda_{\min}^{\star 8/3})$ ,  $1 - (\lambda_{\min}^{\star 4/3}/5)\eta \le \rho < 1$ , and  $C_7$  and  $C_8$  are sufficiently large.

Lemma 6. Suppose that

$$p \ge c_0 \frac{\mu^3 r^2 \log^3 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^*} \le c_1 \frac{\sqrt{p}}{\mu^{3/2} r \sqrt{d \log d}}, \quad \text{and} \quad r \le c_2 \sqrt{\frac{d}{\mu}}$$

for some sufficiently large constant  $c_0 > 0$  and some sufficiently small constant  $c_1, c_2 > 0$ . Assume that the hypotheses (41) hold for the tth iteration and  $\mathcal{E}_{local} \leq c_3/(\mu^{3/2}r)$  for some sufficiently small constant  $c_3 > 0$ . Then with probability at least  $1 - O(d^{-10})$ , one has

$$\|\boldsymbol{U}^{t+1} - \boldsymbol{U}^{\star}\|_{2,\infty} \le \left(C_3 \rho^{t+1} \mathcal{E}_{\text{local}} + C_4 \frac{\sigma}{\lambda_{\min}^{\star}} \sqrt{\frac{d \log d}{p}}\right) \|\boldsymbol{U}^{\star}\|_{2,\infty},$$
(47)

provided that  $0 < \eta \le \lambda_{\min}^{\star 4/3}/(32\lambda_{\min}^{\star 8/3})$ ,  $1 - (\lambda_{\min}^{\star 4/3}/5)\eta \le \rho < 1$ , and  $C_3/(C_5 + C_7)$  and  $C_4/(C_6 + C_8)$  are both sufficiently large.

# 5.2. Analysis for Initialization: Part 1 (Subspace Estimation)

**5.2.1. Key Results.** The aim of this subsection is to demonstrate that the subspace estimate U computed by Algorithm 2 is sufficiently close to the space spanned by the true tensor factors. Given that the columns of  $U^* = [u_1^*, \dots, u_r^*]$  are in general not orthogonal to each other, we shall define  $U_{\text{orth}}^* \in \mathbb{R}^{d \times r}$  as follows (obtained by proper orthonormalization):

$$\boldsymbol{U}_{\text{orth}}^{\star} := \boldsymbol{U}^{\star} (\boldsymbol{U}^{\star \top} \boldsymbol{U}^{\star})^{-\frac{1}{2}}. \tag{48}$$

This matrix  $U_{\text{orth}}^{\star}$  reflects the rank-r principal subspace of  $A^{\star}A^{\star \top} = \sum_{i} \|u_{i}^{\star}\|_{2}^{4} u_{i}^{\star} u_{i}^{\star \top}$ , where we recall that  $A^{\star} \in \mathbb{R}^{d \times d^{2}}$  is the mode-1 matricization of  $T^{\star}$ . In addition, we define the rotation matrix

$$R := \underset{Q \in \mathcal{O}^{r \times r}}{\min} \| UQ - U_{\text{orth}}^{\star} \|_{F}, \tag{49}$$

where  $\mathcal{O}^{r \times r}$  stands for the set of  $r \times r$  orthonormal matrices. This can be viewed as the global rotation matrix that best aligns the two subspaces represented by U and  $U^*_{\text{orth}}$ , respectively.

Equipped with the aforementioned notation, we can invoke (Cai et al. 2021, corollary 1) to arrive at the following lemma, which upper-bounds the distance between our subspace estimate  $\boldsymbol{U}$  and the ground truth  $\boldsymbol{U}_{\text{orth}}^{\star}$ .

**Lemma 7.** There exist some universal constants  $c_0$ ,  $c_1$ ,  $c_2 > 0$  such that if

$$p \ge c_0 \frac{\mu^2 r \log^2 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^*} \le c_1 \frac{\sqrt{p}}{d^{3/4} \sqrt{\log d}}, \quad \text{and} \quad r \le c_2 \sqrt{\frac{d}{\mu'}},$$

then with probability  $1 - O(d^{-10})$ , the subspace estimate **U** computed by Algorithm 2 obeys

$$\| \mathbf{U}\mathbf{R} - \mathbf{U}_{\text{orth}}^{\star} \| \lesssim \mathcal{E}_{\text{se}},$$
 (50a)

$$\| \mathbf{U}\mathbf{R} - \mathbf{U}_{\text{orth}}^{\star} \|_{2,\infty} \lesssim \mathcal{E}_{\text{se}} \sqrt{\frac{\mu r}{d}},$$
 (50b)

where  $\mathbf{U}_{\text{orth}}^{\star}$  and  $\mathbf{R}$  are defined, respectively, in (48) and (49), and

$$\mathcal{E}_{se} := \frac{\mu^2 r \log d}{d^{3/2} p} + \sqrt{\frac{\mu^2 r \log d}{d^2 p}} + \frac{\sigma^2}{\lambda_{\min}^{\star 2}} \frac{d^{3/2} \log d}{p} + \frac{\sigma}{\lambda_{\min}^{\star}} \sqrt{\frac{d \log d}{p}} + \frac{\mu r}{d}.$$
 (51)

In a nutshell, Lemma 7 asserts that, under our sample size, noise, and rank conditions, Algorithm 2 produces reliable estimates of the subspace spanned by the low-rank tensor factors  $\{u_i^\star\}_{1 \leq i \leq r}$ . The theorem quantifies the subspace distance in terms of both the spectral norm and  $\|\cdot\|_{2,\infty}$ , where the latter bound often reflects a considerably stronger sense of proximity compared with the former one.

As it turns out, in order to facilitate analysis for the subsequent stages, we need to introduce a certain leave-one-out sequences as well, which we detail in the next subsection.

**5.2.2.** Leave-One-Out Sequences for Subspace Estimation. The key idea of the leave-one-out analysis is to create auxiliary leave-one-out sequences that are (1) independent of a small fraction of the data and (2) sufficiently close to the true estimates. We introduce the following auxiliary tensor and  $d \times d^2$ -dimensional matrix for each  $1 \le m \le d$ :

$$T^{(m)} := \mathcal{P}_{\Omega_{-m}}(T) + p\mathcal{P}_m(T^*) \in \mathbb{R}^{d \times d \times d}, \tag{52}$$

$$A^{(m)} := \text{mode-1 matricization of } \frac{1}{p} T^{(m)}.$$
 (53)

By construction,  $T^{(m)}$  and  $A^{(m)}$  are independent of  $\mathcal{P}_{\Omega_m}(E)$ , where we recall that

$$\Omega_{-m} := \{(i, j, k) \in \Omega : i \neq m \text{ and } j \neq m \text{ and } k \neq m\},$$
(54)

$$\Omega_m := \{(i, j, k) \in \Omega : i = m \text{ or } j = m \text{ or } k = m\}.$$
(55)

We are now ready to introduce the auxiliary leaveone-out procedure for subspace estimation. Similar to the matrix B in Algorithm 2 (whose eigenspace serves as an estimate of the column space of  $U^*$ ), we define an auxiliary matrix  $B^{(m)} \in \mathbb{R}^{d \times d}$  as follows:

$$\boldsymbol{B}^{(m)} = \mathcal{P}_{\mathsf{off-diag}} (\boldsymbol{A}^{(m)} \boldsymbol{A}^{(m)\top}), \tag{56}$$

where  $\mathcal{P}_{\text{off-diag}}(\cdot)$  (as already defined in Section 3.1) extracts out off-diagonal entries from a matrix. The rationale is simple: it can be easily verified that

$$\mathbb{E}[\boldsymbol{B}^{(m)}] = \boldsymbol{B}^{\star} - \mathcal{P}_{\mathsf{diag}}(\boldsymbol{B}^{\star}), \quad \boldsymbol{B}^{\star} := \boldsymbol{A}^{\star} \boldsymbol{A}^{\star \top}, \tag{57}$$

where  $\mathcal{P}_{\mathsf{diag}}(\cdot)$  extracts out the diagonal entries of the matrix. This gives hope that the eigenspace of  $\mathbf{B}^{(m)}$  is also a reliable estimate of the column space of  $\mathbf{U}^{\star}$ , provided that the diagonal entries of  $\mathbf{B}^{\star}$  are sufficiently small. Consequently, we shall compute  $\mathbf{U}^{0,(m)} \in \mathbb{R}^{d \times r}$  a matrix whose columns are the top-r leading eigenvectors of  $\mathbf{B}^{(m)}$ . The procedure is summarized in Algorithm 5.

**Algorithm 5** (The *m* th Leave-One-Out Sequence for Spectral Initialization)

Spectral Initialization)
1: Let  $\boldsymbol{U}^{(m)}\boldsymbol{\Lambda}^{(m)}\boldsymbol{U}^{(m)\top}$  be the rank-r eigen-decomposition of  $\boldsymbol{B}^{(m)}$  defined in (56).

2: Generate the initial estimate  $\mathbf{U}^{0,(m)} \in \mathbb{R}^{d \times r}$  from  $\mathbf{U}^{(m)} \in \mathbb{R}^{d \times r}$  using Algorithm 6.

The following lemma plays a crucial role in our analysis, which formalizes the fact that the leave-one-out version  $U^{(m)}$  obtained by Algorithm 5 is extremely close to U.

**Lemma 8.** There exist some universal constants  $c_0$ ,  $c_1$ ,  $c_2 > 0$  such that if

$$p \ge c_0 \frac{\mu^2 r \log^2 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^*} \le c_1 \frac{\sqrt{p}}{d^{3/4} \sqrt{\log d}}, \quad \text{and} \quad r \le c_2 \sqrt{\frac{d}{\mu'}},$$

then with probability  $1 - O(d^{-10})$ , the subspace estimate  $U^{(m)}$  computed by Algorithm 5 obeys

$$\|\boldsymbol{U}\boldsymbol{U}^{\top} - \boldsymbol{U}^{(m)}\boldsymbol{U}^{(m)\top}\|_{F} \lesssim \mathcal{E}_{loo}\sqrt{\frac{\mu r}{d}}$$
 (58)

simultaneously for all  $1 \le m \le d$ , where

$$\mathcal{E}_{loo} := \frac{\mu^2 r \log d}{d^{3/2} p} + \sqrt{\frac{\mu^2 r \log d}{d^2 p}} + \frac{\sigma^2}{\lambda_{\min}^{\star 2}} \frac{d^{3/2} \log d}{p} + \frac{\sigma}{\lambda_{\min}^{\star \infty}} \sqrt{\frac{d \log d}{p}}.$$

$$(59)$$

Lemma 8 follows immediately from the analysis of (Cai et al. 2021, lemma 4). As a remark, the construction of the leave-one-out sequences herein is slightly different from the one in (Cai et al. 2021). However, it is straightforward to adapt the proof of (Cai et al. 2021) to the case considered herein. We therefore omit the proof for the sake of brevity.

# 5.3. Analysis for Initialization: Part 2 (Retrieval of Individual Tensor Factors)

**5.3.1. Main Results and Leave-One-Out Sequences.** This section justifies that the procedure presented in Algorithm 3 allows to disentangle the tensor factors. For notational simplicity, we let

$$\overline{u}_{i}^{\star} := u_{i}^{\star} / ||u_{i}^{\star}||_{2}, \quad \lambda_{i}^{\star} := ||u_{i}^{\star}||_{2}^{3}, \quad 1 \leq i \leq d.$$
 (60)

Our result is the following.

**Theorem 3.** Fix any arbitrary small constant  $\delta > 0$ . Assume that

$$p \geq c_0 \frac{\mu^2 r^4 \log^4 d}{d^{3/2}}, \quad \frac{\sigma}{\lambda_{\min}^*} \leq c_1 \frac{\sqrt{p}}{r^{3/2} d^{3/4} \log^2 d'}$$

$$r \leq c_2 \left(\frac{d}{\mu^6 \log^6 d}\right)^{1/6}, \qquad (61)$$

$$L = c_3 r^{2\kappa^2} \log^{3/2} r,$$

$$\epsilon_{th} = c_4 \left\{\frac{\mu r \log d}{d\sqrt{p}} + \frac{\sigma}{\lambda_{\min}^*} \sqrt{\frac{r d \log^2 d}{p}} + \sqrt{\frac{\mu r \log d}{d}}\right\}$$

for some sufficiently large universal constant  $c_0, c_3 > 0$  and some sufficiently small universal constants  $c_1, c_2, c_4 > 0$ . Then, with probability exceeding  $1 - \delta$ , there exists a permutation  $\pi(\cdot): [d] \mapsto [d]$  such that for all  $1 \le i \le r$ , the tensor factors  $\{w^i\}_{i=1}^r$  returned by Algorithm 3 satisfy

$$\| \boldsymbol{w}^{i} - \overline{\boldsymbol{u}}_{\pi(i)}^{\star} \|_{2} \lesssim \frac{\mu r \log d}{d\sqrt{p}} + \frac{\sigma}{\lambda_{\min}^{\star}} \sqrt{\frac{r d \log^{2} d}{p}} + \sqrt{\frac{\mu r \log d}{d}};$$
(62a)

$$\|\boldsymbol{w}^{i} - \overline{\boldsymbol{u}}_{\pi(i)}^{\star}\|_{\infty}$$

$$\leq \left\{ \frac{\mu^{2} \operatorname{rlog}^{4} d}{d^{3/2} p} + \frac{\mu r \log^{3} d}{d \sqrt{p}} + \frac{\sigma^{2}}{\lambda_{\min}^{\star 2}} \frac{d^{3/2} \log^{4} d}{p} + \frac{\sigma}{\lambda_{\min}^{\star 2}} \sqrt{\frac{r d \log^{6} d}{p}} + \sqrt{\frac{\mu r \log^{2} d}{d}} \right\} \sqrt{\frac{\mu r}{d}}; \qquad (62b)$$

$$|\lambda_{i} - \lambda_{\pi(i)}^{\star}| \leq \left\{ \frac{\mu r \log d}{d \sqrt{p}} + \frac{\sigma}{\lambda_{\infty}^{\star}} \sqrt{\frac{r d \log^{2} d}{p}} + \sqrt{\frac{\mu r \log d}{d}} \right\} \lambda_{\pi(i)}^{\star}.$$

In short, this theorem asserts that the estimates returned by Algorithm 3 are—up to global permutation—reasonably close to the ground truth under our

(62c)

sample size and noise conditions. In order to establish this theorem, and in order to provide initial guesses for the leave-one-out GD sequences, we need to produce a leave-one-out sequence tailored to this part of the algorithm. Such auxiliary sequences are generated in a similar spirit as the previous ones, and we summarize them in Algorithm 6. As usual, the resulting leave-one-out estimates  $\{\lambda_i^{(m)}, w^{i,(m)}\}_{i=1}^r$  are statistically independent of  $\mathcal{P}_{\Omega_m}(E)$ .

In what follows, we gather a few key properties of the leave-one-out estimates, which play a crucial role in the analysis.

**Algorithm 6** (The *m* th Leave-One-Out Sequence for Retrieving Individual Tensor Components)

- 1: **Input:** restart number L, threshold  $\epsilon_{th}$ , subspace estimate  $\boldsymbol{U}^{(m)}$  given by Algorithm 5.
- 2: **for**  $\tau = 1, ..., L$  **do**
- 3: Recall the Gaussian vector  $\mathbf{g}^{\tau} \sim \mathcal{N}(0, \mathbf{I}_d)$  generated in Algorithm 3.
- 4:  $(\boldsymbol{\nu}^{\tau,(m)}, \lambda_{\tau}^{(m)}, \operatorname{spec-gap}_{\tau}^{(m)}) \leftarrow \operatorname{Retrieve-one-tensor-factor}(\boldsymbol{T}^{(m)}, p, \boldsymbol{U}^{(m)}, \boldsymbol{g}^{\tau}).$
- 5: Generate tensor factor estimates

$$\begin{split} \{(\boldsymbol{w}^{1,(m)}, \boldsymbol{\lambda}_1^{(m)}), \dots, (\boldsymbol{w}^{r,(m)}, \boldsymbol{\lambda}_r^{(m)})\} \\ &\leftarrow \mathsf{PRUNE}(\{(\boldsymbol{\nu}^{\tau,(m)}, \boldsymbol{\lambda}_\tau^{(m)}.\mathsf{spec-gap}_\tau^{(m)})\}_{\tau=1}^L, \epsilon_{\mathsf{th}}). \end{split}$$

6: **Output:** an initial estimate  $\boldsymbol{U}^{0,(m)} = [(\lambda_1^{(m)})^{1/3} \boldsymbol{w}^{1,(m)}, \dots, (\lambda_r^{(m)})^{1/3} \boldsymbol{w}^{r,(m)}].$ 

**Theorem 4.** Fix any arbitrarily small constant  $\delta > 0$ . Instate the assumptions in Theorem 3. With probability exceeding  $1 - \delta$ , the permutation function stated in Theorem 3 obeys that, for all  $1 \le i \le r$  and all  $1 \le m \le d$ ,

$$\|\boldsymbol{w}^{i} - \boldsymbol{w}^{i,(m)}\|_{2} \lesssim \left\{ \frac{\mu^{2} r \log^{3/2} d}{d^{3/2} p} + \frac{\mu \sqrt{r} \log d}{d \sqrt{p}} + \frac{\sigma^{2}}{\lambda_{\min}^{\star 2}} \frac{d^{3/2} \log^{3/2} d}{p} + \frac{\sigma}{\lambda_{\min}^{\star}} \sqrt{\frac{d \log^{2} d}{p}} \right\} \sqrt{\frac{\mu r}{d}},$$
(63a)

$$\begin{aligned} |\lambda_{i} - \lambda_{i}^{(m)}| & \leq \left\{ \frac{\mu^{2} r \log^{3/2} d}{d^{3/2} p} + \frac{\mu \sqrt{r} \log d}{d \sqrt{p}} + \frac{\sigma^{2}}{\lambda_{\min}^{\star 2}} \frac{d^{3/2} \log^{3/2} d}{p} \right. \\ & + \frac{\sigma}{\lambda_{\min}^{\star}} \sqrt{\frac{d \log^{2} d}{p}} \right\} \sqrt{\frac{\mu r}{d}} \lambda_{\max}^{\star}; \end{aligned}$$

$$(63b)$$

$$\left| \left( \boldsymbol{w}^{i,(m)} - \overline{\boldsymbol{u}}_{\pi(i)}^{\star} \right)_{m} \right| \lesssim \left\{ \frac{\sqrt{\mu r} \log^{7/2} d}{d^{3/2} p} + \frac{\mu r \log^{3} d}{d \sqrt{p}} + \frac{\sigma}{\lambda_{\min}^{\star}} \frac{\log^{4} d}{p} + \frac{\sigma}{\lambda_{\min}^{\star}} \sqrt{\frac{r d \log^{6} d}{p}} + \sqrt{\frac{\mu r \log^{2} d}{d}} \right\} \sqrt{\frac{\mu r}{d}}.$$
(63c)

With Theorems 3–4 in place, we can immediately establish a few desired properties (particularly those specified in Section 5.1) of our initial estimate, as asserted in the following corollary.

**Corollary 3.** Fix any arbitrarily small constant  $\delta > 0$ . Instate the assumptions in Theorem 2. With probability exceeding  $1 - \delta$ , the estimates  $\mathbf{U}^0$  and  $\mathbf{U}^{0,(m)}$  returned by Algorithm 3 and Algorithm 6, respectively, satisfy the hypotheses (41) for t = 0.

## 6. Discussion

The current paper uncovers the possibility of efficiently and stably completing a low-CP-rank tensor from partial and noisy entries. Perhaps somewhat unexpectedly, despite the high degree of nonconvexity, this problem can be solved to optimal statistical accuracy within nearly linear time, provided that the tensor of interest is well conditioned, incoherent, and of constant rank. To the best of our knowledge, this intriguing message has not been shown in the prior literature.

Moving forward, one pressing issue is to understand how to improve the algorithmic and theoretical dependency upon the tensor rank r of the proposed method. Ideally, one would desire a fast algorithm whose sample complexity scales as  $rd^{1.5}$ , an order that is provably achievable by the sum-of-squares hierarchy. Additionally, in contrast to the matrix counterpart where the rank is upper bounded by the matrix dimension, the tensor CP rank is allowed to rise above d, which is commonly referred to as the over-complete case. Unfortunately, our current initialization scheme (i.e., the spectral method) fails to work unless r < d, and our local analysis for GD falls of accommodating the scenario with r > d. It would be of great interest to develop more powerful algorithms-in addition to more refined analysis—to tackle such an important over-complete regime.

Another tantalizing research direction is the exploration of landscape design for tensor completion. As our heuristic discussions as well as other prior work (e.g., Richard and Montanari 2014) suggest, randomly initialized gradient descent tailored to (4) seems unlikely to work, unless the sample size is significantly larger than the computational limit. This might mean either that there exist spurious local minima in the natural nonconvex least-squares formulation (4), or that the optimization landscape of (4) is too flat around some saddle points and hence not amenable to fast computation. It would be interesting to investigate what families of loss functions allow us to rule out bad local minima and eliminate the need of careful initialization, which might be better suited for tensor recovery problems.

Finally, in statistical inference and decision making, one might not be simply satisfied with obtaining a reliable estimate for each missing entry, but would also like to report a short confidence interval which is likely to contain the true entry. This boils down to the fundamental task of uncertainty quantification for tensor completion, which we leave to future investigation.

### **Acknowledgments**

The authors thank Lanqing Yu for many helpful discussions and Yuling Yan for proofreading the paper. This paper was first submitted in November 2019 and was accepted in part to NeurIPS 2019 (Cai et al. 2019).

### **Endnotes**

<sup>1</sup> We focus on symmetric order-3 tensors primarily for simplicity of presentation. Many of our findings naturally extend to the more general case with asymmetric tensors of possibly higher order. Detailed discussions are deferred to Section EC.7 in the e-companion due to the space limits.

<sup>2</sup> Here, a tensor  $T \in \mathbb{R}^{d \times d \times d}$  is said to be *symmetric* if  $T_{j,k,l} = T_{k,j,l} = T_{k,l,j} = T_{l,k,l} = T_{l,l,k} = T_{l,l,k}$  for all  $1 \le j,k,l \le d$ .

#### References

Abbe E, Fan J, Wang K, Zhong Y (2017) Entrywise eigenvector analysis of random matrices with low expected rank. *Ann. Statist*. 48(3):1452–1474.

Anandkumar A, Ge R, Janzamin M (2014b) Guaranteed non-orthogonal tensor decomposition via alternating rank-1 updates. Preprint, submitted February 21, https://arxiv.org/abs/1402.5180.

Anandkumar A, Ge R, Hsu D, Kakade SM, Telgarsky M (2014a) Tensor decompositions for learning latent variable models. *J. Machine Learn. Res.* 15(1):2773–2832.

Ashraphijuo M, Wang X (2017) Fundamental conditions for low-CP-rank tensor completion. *J. Machine Learn. Res.* 18(1):2116–2145.

Barak B, Moitra A (2016) Noisy tensor completion via the sum-of-squares hierarchy. *Proc. 29th Annual Conf. Learning Theory, New York, June 23–26*, 417–445.

Bubeck S (2015) Convex optimization: Algorithms and complexity. *Foundations Trends Machine Learning* 8(3–4):231–357.

Cai C, Li G, Poor HV, Chen Y (2019) Nonconvex low-rank tensor completion from noisy data. Wallach H, Larochelle H, Beygelzimer A, d'Alché-Buc F, Fox E, Garnett R, eds. Adv. Neural Inform. Processing Systems 32 (Curran Associates, Red Hook, NY), 1863–1874.

Cai C, Li G, Chi Y, Poor HV, Chen Y (2021) Subspace estimation from unbalanced and incomplete data matrices:  $\ell_{2,\infty}$  statistical guarantees. *Ann. Statist.* 49(2):944–967.

Candès E, Recht B (2009) Exact matrix completion via convex optimization. Foundations Comput. Math. 9(6):717–772.

Candès EJ, Li X, Soltanolkotabi M (2015) Phase retrieval via Wirtinger flow: Theory and algorithms. IEEE Trans. Inform. Theory 61(4):1985–2007.

Chen Y, Candès EJ (2017) Solving random quadratic systems of equations is nearly as easy as solving linear systems. *Comm. Pure Appl. Math.* 70(5):822–883.

Chen Y, Chi Y (2014) Robust spectral compressed sensing via structured matrix completion. *IEEE Trans. Inform. Theory* 60(10): 6576–6601.

Chen Y, Chi Y (2018) Harnessing structures in big data via guaranteed low-rank matrix estimation: Recent theory and fast algorithms via convex and nonconvex optimization. *IEEE Signal Processing Magazine* 35(4):14–31.

- Chen Y, Wainwright MJ (2015) Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees. Preprint, submitted September 10, https://arxiv.org/abs/1509.03025.
- Chen J, Liu D, Li X (2019a) Nonconvex rectangular matrix completion via gradient descent without  $\ell_{2,\infty}$  regularization. Preprint, submitted January 18, https://arxiv.org/abs/1901.06116v1.
- Chen H, Raskutti G, Yuan M (2019b) Non-convex projected gradient descent for generalized low-rank tensor regression. J. Machine Learn. Res. 20(1):172–208.
- Chen Y, Chi Y, Fan J, Ma C (2019c) Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval. *Math. Programming* 176(1–2):5–37.
- Chen Y, Fan J, Ma C, Wang K (2019d) Spectral method and regularized MLE are both optimal for top-K ranking. *Ann. Statist.* 47 (4):2204–2235.
- Chen Y, Fan J, Ma C, Yan Y (2019e) Inference and uncertainty quantification for noisy matrix completion. *Proc. Natl. Acad. Sci. USA* 116(46):22931–22937.
- Chen Y, Chi Y, Fan J, Ma C, Yan Y (2021) Noisy matrix completion: Understanding statistical guarantees for convex relaxation via nonconvex optimization. *SIAM J. Optim.* 30(4):3098–3121.
- Cheng JY, Zhang T, Alley MT, Uecker M, Lustig M, Pauly JM, Vasanawala SS (2017) Comprehensive multi-dimensional MRI for the simultaneous assessment of cardiopulmonary anatomy and physiology. *Sci. Rep.* 7(1):5330.
- Chi Y, Lu Y, Chen Y (2019) Nonconvex optimization meets lowrank matrix factorization: An overview. *IEEE Trans. Signal Proc*essing 67(20):5239–5269.
- Davenport MA, Romberg J (2016) An overview of low-rank matrix recovery from incomplete observations. *IEEE J. Selected Topics Signal Processing* 10(4):608–622.
- Ding L, Chen Y (2018) Leave-one-out approach for matrix completion: Primal and dual analysis. Preprint, submitted March 20, https://arxiv.org/abs/1803.07554.
- Driggs D, Becker S, Boyd-Graber J (2019) Tensor robust principal component analysis: Better recovery with atomic norm regularization. Preprint, submitted January 30, https://arxiv.org/abs/ 1901.10991.
- El Karoui N (2015) On the impact of predictor geometry on the performance on high-dimensional ridge-regularized generalized robust regression estimators. *Probab. Theory Related Fields* 170: 05-175
- Ely G, Aeron S, Hao N, Kilmer ME (2013) 5D and 4D pre-stack seismic data completion using tensor nuclear norm (TNN). *SEG Tech. Program Expanded Abstr.* (Society of Exploration Geophysicists, Tulsa, OK), 3639–3644.
- Farias VF, Li AA (2019) Learning preferences with side information. *Management Sci.* 65(7):3131–3149.
- Gandy S, Recht B, Yamada I (2011) Tensor completion and low-nrank tensor recovery via convex optimization. *Inverse Problems* 27(2):025010.
- Ge R, Ma T (2017) On the optimization landscape of tensor decompositions. von Luxburg U, Guyon I, Bengio S, Wallach H, Fergus R, eds. *Adv. Neural Inform. Processing Systems* 31 (Curran Associates, Red Hook, NY), 3653–3663.
- Ge R, Huang F, Jin C, Yuan Y (2015) Escaping from saddle points: Online stochastic gradient for tensor decomposition. *Proc. 28th Conf. Learning Theory, Paris*, July 3–6, 797–842.
- Gilboa D, Buchanan S, Wright J (2018) Efficient dictionary learning with gradient descent. Preprint, submitted September 27, https://arxiv.org/abs/1809.10313.
- Goldfarb D, Qin Z (2014) Robust low-rank tensor recovery: Models and algorithms. SIAM J. Matrix Anal. Appl. 35(1):225–253.
- Gross D (2011) Recovering low-rank matrices from few coefficients in any basis. *IEEE Trans. Inform. Theory* 57(3):1548–1566.

- Hao B, Zhang A, Cheng G (2020) Sparse and low-rank tensor estimation via cubic sketchings. *IEEE Trans. Inform. Theory* 66(9): 5927–5964.
- Hao B, Wang B, Wang P, Zhang J, Yang J, Sun WW (2019) Sparse tensor additive regression. Submitted March 31, https://arxiv.org/abs/1904.00479.
- Hillar CJ, Lim LH (2013) Most tensor problems are np-hard. *J. ACM* 60(6):45.
- Hopkins SB, Schramm T, Shi J, Steurer D (2016) Fast spectral algorithms from sum-of-squares proofs: tensor decomposition and planted sparse vectors. *Proc.48th Annual ACM Sympos. Theory Comput.* (ACM, New York), 178–191.
- Huang B, Mu C, Goldfarb D, Wright J (2015) Provable models for robust low-rank tensor completion. *Pacific J. Optim.* 11(2): 339–364.
- Jain P, Oh S (2014) Provable tensor factorization with missing data. Ghahramani Z, Welling M, Cortes C, Lawrence N, Weinberger KQ, eds. Adv. Neural Inform. Processing Systems 27 (Curran Associates, Red Hook, NY), 1431–1439.
- Ji TY, Huang TZ, Zhao XL, Ma TH, Liu G (2016) Tensor completion using total variation and low-rank matrix factorization. *Inform. Sci.* 326:243–257.
- Kasai H, Mishra B (2016) Low-rank tensor completion: a riemannian manifold preconditioning approach. Balcan MF, Weinberger KQ, eds. Proc. 33rd Internat. Conf. Machine Learning, New York, June 20–22, 1012–1021.
- Keshavan RH, Montanari A, Oh S (2010a) Matrix completion from a few entries. *IEEE Trans. Inform. Theory* 56(6):2980–2998.
- Keshavan RH, Montanari A, Oh S (2010b) Matrix completion from noisy entries. *J. Machine Learn. Res.* 11(69):2057–2078.
- Kilmer ME, Braman K, Hao N, Hoover RC (2013) Third-order tensors as operators on matrices: A theoretical and computational framework with applications in imaging. *SIAM J. Matrix Anal. Appl.* 34(1):148–172.
- Kim HJ, Ollila E, Koivunen V, Croux C (2013) Robust and sparse estimation of tensor decompositions. *Proc.* 2013 IEEE Global Conf. Signal Inform. Processing (IEEE, Piscataway, NJ), 965–968.
- Kolda TG (2001) Orthogonal tensor decompositions. SIAM J. Matrix Anal. Appl. 23(1):243–255.
- Kolda TG, Bader BW (2009) Tensor decompositions and applications. SIAM Rev. 51(3):455–500.
- Kreimer N, Stanton A, Sacchi MD (2013) Tensor completion based on nuclear norm minimization for 5d seismic data reconstruction. *Geophysics* 78(6):273–284.
- Krishnamurthy A, Singh A (2013) Low-rank matrix and tensor completion via adaptive sampling. Burges CJC, Bottou L, Welling M, Ghahramani Z, Weinberger KQ, eds. *Adv. Neural Inform. Processing Systems* 26 (Curran Associates, Red Hook, NY), 836–844.
- Li Q, Tang G (2017) Convex and nonconvex geometries of symmetric tensor factorization. *Proc. 51st Asilomar Conf. Signals Systems Comput.* (IEEE, Piscataway, NJ), 305–309.
- Li X, Ye Y, Xu X (2017) Low-rank tensor completion with total variation for visual data inpainting. Singh S, Markovitch S, eds. *Proc. 31st AAAI Conf. Artificial Intelligence* (AAAI Press, San Francisco), 2210–2216.
- Li Q, Zhu Z, Tang G (2019) The non-convex geometry of low-rank matrix optimization. *Inform. Inference* 8(1):51–96.
- Liu J, Musialski P, Wonka P, Ye J (2013) Tensor completion for estimating missing values in visual data. *IEEE Trans. Pattern Anal. Machine Intelligence* 35(1):208–220.
- Liu Y, Shang F, Cheng H, Cheng J, Tong H (2014) Factor matrix trace norm minimization for low-rank tensor completion. Zaki M, Obradovic Z, Tan PN, Banerjee A, Kamath C, Parthasarathy S, eds. Proc. 2014 SIAM Internat. Conf. Data Mining (SIAM, Philadelphia), 866–874.
- Lounici K (2014) High-dimensional covariance matrix estimation with missing observations. *Bernoulli* 20(3):1029–1058.

- Lu C, Feng J, Chen Y, Liu W, Lin Z, Yan S (2016) Tensor robust principal component analysis: Exact recovery of corrupted lowrank tensors via convex optimization. *Proc. IEEE Conf. Comput. Vision Pattern Recognition* (IEEE, Piscataway, NJ), 5249–5257.
- Ma C, Wang K, Chi Y, Chen Y (2017) Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion and blind deconvolution. Foundations Comput. Math. 20:451–463.
- Mišić VV, Perakis G (2020) Data analytics in operations management: A review. *Manufacturing Service Oper. Management* 22(1): 158–169.
- Montanari A, Sun N (2018) Spectral algorithms for tensor completion. *Comm. Pure Appl. Math.* 71(11):2381–2425.
- Mu C, Huang B, Wright J, Goldfarb D (2014) Square deal: Lower bounds and improved relaxations for tensor recovery. Xing EP, Jebara T, eds. Proc. 31st Internat. Conf. Machine Learning, Beijing, June 22–24, 73–81.
- Pananjady A, Wainwright MJ (2019) Instance-dependent  $\ell_{\infty}$ -bounds for policy evaluation. Preprint, submitted September 19, https://arxiv.org/abs/1909.08749.
- Pawlowski C (2019) Machine learning for problems with missing and uncertain data with applications to personalized medicine. Unpublished doctoral dissertation, Massachusetts Institute of Technology.
- Potechin A, Steurer D (2017) Exact tensor completion with sum-of-squares. Kale S, Shamir O, eds. *Proc. 2017 Conf. Learning Theory, Amsterdam*, July 7–10, 1619–1673.
- Rauhut H, Schneider R, Stojanac Ž (2017) Low rank tensor recovery via iterative hard thresholding. *Linear Algebra Appl.* 523:220–262.
- Richard E, Montanari A (2014) A statistical model for tensor PCA. Ghahramani Z, Welling M, Cortes C, Lawrence ND, Weinberger KQ, eds. *Adv. Neural Inform. Processing, Systems* 27 (Curran Associates, Red Hook, NY), 2897–2905.
- Romera-Paredes B, Pontil M (2013) A new convex relaxation for tensor completion. Burges CJC, Bottou L, Welling M, Ghahramani Z, Weinberger KQ, eds. *Adv. Neural Inform. Processing Systems* 26 (Curran Associates, Red Hook, NY), 2967–2975.
- Semerci O, Hao N, Kilmer ME, Miller EL (2014) Tensor-based formulation and nuclear norm regularization for multienergy computed tomography. *IEEE Trans. Image Processing* 23(4): 1678–1693.
- Shah D, Yu CL (2019) Iterative collaborative filtering for sparse noisy tensor estimation. Preprint, submitted August 3, https:// arxiv.org/abs/1908.01241.
- Sidiropoulos ND, De Lathauwer L, Fu X, Huang K, Papalexakis EE, Faloutsos C (2017) Tensor decomposition for signal processing and machine learning. *IEEE Trans. Signal Processing* 65(13): 3551–3582.
- Soroushmehr SR, Najarian K (2016) Transforming big data into computational models for personalized medicine and healthcare. *Dialogues Clinical Neuroscience* 18(3):339–343.
- Steinlechner M (2016) Riemannian optimization for high-dimensional tensor completion. SIAM J. Sci. Comput. 38(5):S461–S484.
- Sun R, Luo ZQ (2016) Guaranteed matrix completion via non-convex factorization. *IEEE Trans. Inform. Theory* 62(11):6535–6579.
- Sun WW, Lu J, Liu H, Cheng G (2017) Provable sparse tensor decomposition. J. Royal Statist. Soc. Ser. B. Statist. Methodology 79(3): 899–916.
- Tan YS, Vershynin R (2019) Online stochastic gradient descent with arbitrary initialization solves non-smooth, non-convex phase retrieval. Preprint, submitted October 28, https://arxiv.org/abs/ 1910.12837.
- Tang G, Shah P (2015) Guaranteed tensor decomposition: A moment approach. Bach F, Blei D, eds. Proc. 32nd Internat. Conf. Machine Learning, Lille, France, July 7–9, 1491–1500.

- Tomioka R, Hayashi K, Kashima H (2010) Estimation of low-rank tensors via convex optimization. Preprint, submitted October 5, https://arxiv.org/abs/1010.0789.
- Wang W, Aggarwal V, Aeron S (2016) Tensor completion by alternating minimization under the tensor train (TT) model. Preprint, submitted September 19, https://arxiv.org/abs/1609.05587.
- Xia D, Yuan M (2017) On polynomial time methods for exact low rank tensor completion. Preprint, submitted February 22, https://arxiv.org/abs/1702.06980.
- Xia D, Yuan M, Zhang CH (2017) Statistically optimal and computationally efficient low rank tensor completion from noisy entries. Preprint, submitted November 14, https://arxiv.org/abs/1711.04934.
- Xu Y, Yin W (2013) A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion. *SIAM J. Imaging Sci.* 6(3): 1758–1789.
- Xu Y, Hao R, Yin W, Su Z (2015) Parallel matrix factorization for low-rank tensor completion. *Inverse Problems Imaging* 9(2):601–624.
- Yao Q (2018) Scalable tensor completion with nonconvex regularization. Preprint, submitted July 23, https://arxiv.org/abs/1807.08725v1.
- Yi X, Park D, Chen Y, Caramanis C (2016) Fast algorithms for robust PCA via gradient descent. Lee D, Sugiyama M, Luxburg U, Guyon I, Garnett R, eds. *Adv. Neural Inform. Processing Systems* 29 (Curran Associates, Red Hook, NY), 4152–4160.
- Ying J, Lu H, Wei Q, Cai JF, Guo D, Wu J, Chen Z, Qu X (2017) Hankel matrix nuclear norm regularized tensor completion for n-dimensional exponential signals. *IEEE Trans. Signal Processing* 65(14):3702–3717.
- Yuan M, Zhang CH (2016) On tensor completion via nuclear norm minimization. *Found. Comput. Math.* 16(4):1031–1068.
- Yuan M, Zhang CH (2017) Incoherent tensor norms and their applications in higher order tensor completion. *IEEE Trans. Inform. Theory* 63(10):6753–6766.
- Zhang A (2019) Cross: Efficient low-rank tensor completion. *Ann. Statist.* 47(2):936–964.
- Zhang Z, Aeron S (2017) Exact tensor completion using t-SVD. *IEEE Trans. Signal Processing* 65(6):1511–1526.
- Zhang A, Xia D (2018) Tensor SVD: Statistical and computational limits. *IEEE Trans. Inform. Theory* 64(11):7311–7338.
- Zhong Y, Boumal N (2018) Near-optimal bound for phase synchronization. SIAM J. Optim. 28(2):989–1016.
- Zoubir AM, Koivunen V, Ollila E, Muma M (2018) Robust Statistics for Signal Processing (Cambridge University Press, Cambridge, UK).
- **Changxiao Cai** is a PhD student in the Department of Electrical Engineering at Princeton University. He received the BE in electronic engineering from Tsinghua University in 2016. His research interests include machine learning, high-dimensional statistics, optimization, and information theory.
- **Gen Li** is a PhD student in the Department of Electronic Engineering at Tsinghua University. He received the BE degree in electronic engineering from Tsinghua University in 2016. His recent research interests include machine learning and nonconvex optimization, high-dimensional statistics, and reinforcement learning.
- H. Vincent Poor is the Michael Henry Strater University Professor of Electrical Engineering at Princeton University. His research interests are in the areas of information theory, machine learning, and network science and their applications in wireless networks, energy systems, and related fields. Recent recognition of his

work includes the 2017 IEEE Alexander Graham Bell Medal and a DEng honoris causa from the University of Waterloo awarded in 2019.

Y. Chen is an assistant professor in the Department of Electrical Engineering at Princeton University. His research interests include high-dimensional statistics, convex and nonconvex optimization, reinforcement learning, and information theory. His most recent awards include the 2020 Army Research Office Young Investigator Award, the 2020 Princeton Graduate Mentoring Award, and the 2020 International Congress of Chinese Mathematicians Best Paper Award (gold medal).