

---

# Scalable and Sound Low-Rank Tensor Learning

---

Hao Cheng<sup>‡</sup>      Yaoliang Yu<sup>\*</sup>      Xinhua Zhang<sup>\*</sup>      Eric Xing<sup>\*</sup>      Dale Schuurmans<sup>†</sup>  
University of Washington<sup>‡</sup>    CMU<sup>\*</sup>    University of Illinois at Chicago<sup>\*</sup>    University of Alberta<sup>†</sup>

## Abstract

Many real-world data arise naturally as tensors. Equipped with a low rank prior, learning algorithms can benefit from exploiting the rich dependency encoded in a tensor. Despite its prevalence in low-rank matrix learning, trace norm ceases to be tractable in tensors and therefore most existing works resort to matrix unfolding. Although some theoretical guarantees are available, these approaches may lose valuable structure information and are not scalable in general. To address this problem, we propose directly optimizing the tensor trace norm by approximating its *dual* spectral norm, and we show that the approximation bounds can be *efficiently* converted to the original problem via the generalized conditional gradient algorithm. The resulting approach is scalable to large datasets, and matches state-of-the-art recovery guarantees. Experimental results on tensor completion and multitask learning confirm the superiority of the proposed method.

## 1 Introduction

Real-world data are complex and usually exhibit rich structures that learning algorithms can significantly benefit from. For instance, data in the vectorial form can be sparse while in the matrix form can have low rank. More general multi-way correlations are often observed in tensor form data, i.e. multi-dimensional arrays [1]. Examples include multi-channel images, video sequences, chemical compound processes, brain EEG signals, etc. Additionally, tensors are often used as important tools for modeling and estimation. Most notably, they can be applied as the high order moments in latent variable models [2] such as independent com-

ponent analysis, mixture and topic models, etc. Due to the wide applicability, a lot of works have focused on factorizing, completing, and learning a tensor of interest.

The most straightforward solution to low-rank tensor decomposition is arguably alternating least squares (ALS) or block-coordinate based methods [1, 3]. With relatively high efficiency, they are generically applicable in practice. In the case of completing a  $d \times d \times d$  tensor using  $z$  observed entries, the complexity is only  $O(zr)$ , where  $r$  is the intended rank. However, they are based on nonconvex optimization whose globally optimal solution is generally intractable to find. Therefore, the analysis of sample complexity and recovery bounds is generally hard, leaving the recent results restricted to strong assumptions such as the absence of noise or the access to the truth rank [e.g. 2, 4, 5].

A prevalent class of theoretically sound approaches are based on convex relaxations of the rank function. In low-rank matrix learning, the utilization of trace norm has achieved remarkable success since the seminal work of [6]. Extension to tensors is conceptually straightforward, based on the generic framework of atomic norm [7]. Recently, [8] provided excellent theoretical support for the tensor trace norm. However, many computational tractability properties of matrices do not carry over to tensors. For instance, the tensor rank can (far) exceed the maximum dimension, and a low-rank approximation may not even exist [9]. Indeed, it is also NP-hard to decide the tensor rank and tensor trace norm, see e.g. [10].

To work around the tractability barrier, approximations are in order. A large body of algorithms unfold the tensor into matrices, and then apply well established matrix techniques, e.g., minimizing the (matrix) trace norm as a proxy of the rank [11–16]. However, this may be unsatisfying because a genuinely “low-rank” tensor can nevertheless have large or even full rank in *all* matrix unfoldings. It also brings about scalability issues as the most efficient implementation of proximal gradient methods or ADMM costs at least  $O(zd)$  time. As a result, most matricization based algorithms do not scale to large tensors.

---

Appearing in Proceedings of the 19<sup>th</sup> International Conference on Artificial Intelligence and Statistics (AISTATS) 2016, Cadiz, Spain. JMLR: W&CP volume 51. Copyright 2016 by the authors.

The aim of this paper is therefore to design approximate models which are both scalable and theoretically sound. Instead of directly approximating the tensor trace norm, we resort to its dual norm—tensor spectral norm, because it is much easier to approximate. We provide a simple algorithm to compute the tensor spectral norm with sound approximation bounds, and its computational cost is similar to that of ALS (Section 4). A key challenge is then to utilize this approximation as an oracle, and to find a solution for the original optimization problem with suboptimality guarantee. We show that this translation can be achieved by using the recent generalized conditional gradient algorithm [17, 18], and furthermore  $O(1/\epsilon)$  rates of convergence can be derived for this approximate oracle setting. Finally we study the sample complexity of the proposed algorithm in Section 5, and demonstrate that it does match the state-of-the-art bounds of matricization based models. Therefore, we obtain an algorithm which is efficient in both computation and inference. In experiments on tensor completion and multitask learning, our method significantly outperforms existing methods in generalization and efficiency.

## 2 Preliminaries

We first review some preliminaries. A tensor is identified as a multi-dimensional array  $\mathcal{A} \in \mathbb{R}^{d_1 \times d_2 \times \dots \times d_K}$ , where  $K$  is the order and  $d_k$  is the dimension of the  $k$ -th mode. Without loss of generality, we assume throughout that  $d_1 \geq d_2 \geq \dots \geq d_K$ . Following [1], we define the mode- $k$  multiplication of a tensor with some matrix  $U \in \mathbb{R}^{\hat{d}_k \times d_k}$  as

$$(\mathcal{A} \times_k U)_{i_1, \dots, i_{k-1}, \hat{i}_k, i_{k+1}, \dots, i_K} = \sum_{i_k=1}^{d_k} \mathcal{A}_{i_1, \dots, i_{k-1}, i_k, i_{k+1}, \dots, i_K} U_{i_k, \hat{i}_k}.$$

Note that the dimension of the  $k$ -th mode of the product  $\mathcal{A} \times_k U$  changes from  $d_k$  to  $\hat{d}_k$ . As usual, we define the inner product of two tensors with the same size as  $\langle \mathcal{A}, \mathcal{B} \rangle := \sum_{i_1, \dots, i_K} \mathcal{A}_{i_1, \dots, i_K} \mathcal{B}_{i_1, \dots, i_K}$ . and the induced Frobenius norm  $\|\mathcal{A}\|_F := \sqrt{\langle \mathcal{A}, \mathcal{A} \rangle}$ .

We can unfold (or flatten) a tensor into a 2-D matrix as follows. For any  $1 \leq k \leq K$ ,  $\mathcal{A}_{(k)} \in \mathbb{R}^{d_k \times (\prod_{j \neq k} d_j)}$ , with its  $(i_k, 1 + \sum_{j=1, j \neq k}^K (i_j - 1) \prod_{m=j+1, m \neq k}^K d_m)$ -th entry being  $\mathcal{A}_{i_1, \dots, i_K}$ . Like matrices, we can decompose a tensor into a linear combination of *primitives*. First, we define rank-1 (pure) tensors as  $\mathcal{T} = \mathbf{u} \otimes \mathbf{v} \otimes \dots \otimes \mathbf{z}$ , where the  $K$  factors  $\mathbf{u} \in \mathbb{R}^{d_1}, \dots, \mathbf{z} \in \mathbb{R}^{d_K}$  form the outer product  $\mathcal{T}_{i_1, i_2, \dots, i_K} = u_{i_1} v_{i_2} \dots z_{i_K}, \forall i_1, \dots, i_K$ . Then, the Candecomp-Parafac (CP) decomposition of the tensor  $\mathcal{A}$  refers to

$$\mathcal{A} = \sum_{i=1}^r \mathbf{u}_i \otimes \mathbf{v}_i \otimes \dots \otimes \mathbf{z}_i, \quad (1)$$

and the lowest possible value of  $r$  is called the *rank* of  $\mathcal{A}$ . Surprisingly, the tensor rank could well exceed any dimension  $d_k$ . In general, a well known bound states

$$\forall k, \text{rank}(\mathcal{A}_{(k)}) \leq \text{rank}(\mathcal{A}) \leq \prod_{j \neq k} \text{rank}(\mathcal{A}_{(j)}), \quad (2)$$

where recall that  $\text{rank}(\mathcal{A}_{(k)})$  is upper bounded by the  $k$ -th dimension  $d_k$ . Note that even when the tensor rank is small, the mode- $k$  unfolding  $\mathcal{A}_{(k)}$  can still have large or even full matrix rank for all  $k$ , which suggests that treating a genuine tensor as a folded matrix could lose a lot of structure.

Numerically finding a CP decomposition or the tensor rank is intractable [10]. Block coordinate descent, i.e. alternatively optimizing each factor matrix with all others fixed, can lead to a *local* optimum. We caution that when  $K > 2$ , deflation, i.e. finding the rank-1 factors successively, can be far from optimal [19].

## 3 Low-Rank Tensor Learning

Learning a low-rank tensor can be formulated as the following optimization problem:

$$\inf_{\mathcal{W}} \ell(\mathcal{W}) + \lambda \cdot \text{rank}(\mathcal{W}), \quad (3)$$

where  $\lambda \geq 0$  is the regularization constant. The loss function  $\ell$  measures the discrepancy between the estimate tensor  $\mathcal{W}$  and the data tensor  $\mathcal{X}$  (suppressed in our notation). Some useful instantiations of  $\ell$  can be found in Section 6, and here we use as a concrete example the tensor completion problem, with  $\ell(\mathcal{W}) = \ell_{\text{tc}}(\mathcal{W}) = \|\mathcal{P}(\mathcal{W} - \mathcal{X})\|_F^2$ , where  $\mathcal{P} : \mathbb{R}^{d_1 \times \dots \times d_K} \rightarrow \mathbb{R}^{d_1 \times \dots \times d_K}$  is the (linear) sampling operator that fills unobserved entries with 0. We are particularly interested in finding a low-rank estimate, induced by penalizing the tensor rank. This bias can be beneficial in multiple ways: a). It is a valid prior in many applications; b). It leads to simpler models with potentially better generalization performance; c). It improves scalability by reducing the cost in computation and storage, as we shall see.

Although appealing in theory, the tensor rank is computationally intractable [10]. For matrices, the trace norm relaxation [6, 20] has been remarkably successful, serving as a convex surrogate of the matrix rank. It is thus natural to extend this principle to tensors, using the atomic norm framework of [7]. We start with a set of atoms (primitives), convexify it if necessary, and then construct the Minkowski gauge function as the appropriate regularizer for promoting the atoms. It is demonstrated in [7] that this yields the “best” convex regularizer in an appropriate sense.

Adapting to the low-rank tensor learning problem (3), we choose the atoms to be rank-1 tensors:

$$\mathbf{A} := \{\mathbf{u}_1 \otimes \dots \otimes \mathbf{u}_K : \forall k, \mathbf{u}_k \in \mathbb{R}^{d_k}, \|\mathbf{u}_k\|_2 \leq 1\}. \quad (4)$$

Upon convexification we can construct the tensor trace norm [TTN, e.g. 21, 22]:

$$\|\mathcal{A}\|_{\text{tr}} := \inf\{\rho \geq 0 : \mathcal{A} \in \rho \cdot \text{conv}(\mathbf{A})\}. \quad (5)$$

Clearly (5) is a norm on  $\mathbb{R}^{d_1 \times \dots \times d_K}$ , and its dual norm—tensor spectral norm (TSN)—is given by

$$\begin{aligned} \|\mathcal{A}\|_{\text{sp}} &:= \max_{\|\mathcal{B}\|_{\text{tr}} \leq 1} \langle \mathcal{A}, \mathcal{B} \rangle = \max_{\mathcal{B} \in \mathbf{A}} \langle \mathcal{A}, \mathcal{B} \rangle \quad (6) \\ &= \max_{\forall k, \|\mathbf{u}_k\|_2 \leq 1} \langle \mathcal{A}, \mathbf{u}_1 \otimes \dots \otimes \mathbf{u}_K \rangle. \quad (7) \end{aligned}$$

Specializing to matrices ( $K = 2$ ) we recover the familiar trace norm and the usual spectral norm.

### 3.1 Main formulation

We can now obtain a convex relaxation of (3) by replacing the tensor rank with the TTN,

$$\min_{\mathcal{W}} \ell(\mathcal{W}) + \lambda \cdot \|\mathcal{W}\|_{\text{tr}}. \quad (8)$$

Assuming  $\ell$  is convex, (8) is a convex problem. This convexity appears to be a significant step towards tractability, and numerous efficient algorithms have been proposed for the matrix case ( $K = 2$ ). However, almost no<sup>1</sup> significant effort has been made for  $K \geq 3$ . [8], for instance, proved some theoretical advantages of using TTN, but no practical algorithm was provided. The underlying reason is that TTN, albeit being convex, is itself intractable [10]! Consequently, one might wonder what is the benefit of replacing the original intractable problem (3) with another convex but still intractable one (8)? The answer in short is: the latter provides more favorable approximation bounds and facilitates more effective algorithms in practice—not all intractable problems are equally “hard”.

Indeed, it is a standard practice to attack NP-hard problems with *approximate* guarantees (additive or multiplicative). We recall the definition of a multiplicative  $\alpha$ -approximate algorithm:

**Definition 1.** Consider a class of optimization problems  $0 \leq \text{OPT}_f = \inf_{\mathbf{w}} f(\mathbf{w})$ , where  $f$  belongs to some function class  $\mathcal{F}$ . Fix  $\alpha > 0$ . An algorithm is  $\alpha$ -approximate if for all  $f \in \mathcal{F}$  it always outputs  $\mathbf{w}_f^*$  such that  $f(\mathbf{w}_f^*) \leq \frac{1}{\alpha} \text{OPT}_f$ .

There is a hierarchy of NP-hard problems that can be solved approximately: some admits polynomial time approximation (e.g. Knapsack); some admits a constant approximation (e.g. max-cut); and still some has guarantees depending on the problem size, which is the case in (6) as will be shown below. We will develop an efficient  $\alpha$ -approximate algorithm for solving the relaxation (8). The convexity in (8), although not leading directly to tractability, will play a key role in our reasoning and development.

<sup>1</sup>After this work was completed, N. Rao kindly brought to our attention a related work [23], which also applied an improved conditional gradient to tensor completion, but no approximation guarantee was provided.

## 4 Efficient Approximate Algorithm

The most straightforward attempt to solve (8) is to approximate the TTN. However, different from the matrix trace norm, it is not the sum of (natural) “singular values”, and the definition in (5) does not provide an explicit form that allows convenient and analyzable approximation. By contrast, the definition of TSN in (6) exhibits much “simpler” and more explicit structure (though also NP-hard), where the constraints decouple over modes with variables  $\mathbf{u}_k$  lying in standard Euclidean balls. So a natural strategy is to first approximate the TSN with good guarantee, and then design optimization algorithms that convert this bound to the optimization bound of the original problem (8).

### 4.1 Approximating the Tensor Spectral Norm

To design an efficient algorithm for approximating TSN, we first resort to a celebrated result from convex geometry [24]. Let  $\mathbf{B}_2^d$  be the Euclidean norm ball in a  $d$  dimensional space ( $d$  is a superscript in  $\mathbf{B}_2^d$ ). We can find in polynomial time (at most)  $d$  points  $\{\mathbf{p}_i\}$ , such that their convex hull  $\mathbf{P}^d$  (a polytope) satisfies  $\mathbf{P}^d \subseteq \mathbf{B}_2^d \subseteq c\sqrt{d}\mathbf{P}^d$ , where  $c$  is some universal constant. In particular we can use the unit  $L_1$  ball  $\mathbf{B}_1^d$ : clearly  $\mathbf{B}_1^d \subseteq \mathbf{B}_2^d \subseteq \sqrt{d}\mathbf{B}_1^d$ . By counting the volume, it can be proved that the factor  $\sqrt{d}$  here is the best possible [25]. Specializing to the TSN, we simply replace each Euclidean ball constraint with an  $L_1$  ball, and evaluate the inner product in (6) at each of the vertices of the polytope. Since the matrix trace norm is tractable, we need only execute this polytopal approximation for the last  $K - 2$  modes, i.e., solving the approximation

$$\max \langle \mathcal{A}, \mathbf{u}_1 \otimes \dots \otimes \mathbf{u}_K \rangle, \text{ s.t. } \forall i \geq 3, \mathbf{u}_i \in \mathbf{P}^{d_i}, \quad (9)$$

where also  $\mathbf{u}_1 \in \mathbf{B}_2^{d_1}, \mathbf{u}_2 \in \mathbf{B}_2^{d_2}$ .

For each of the  $\prod_{k=3}^K d_k$  vertices  $\mathbf{p}^3 \otimes \dots \otimes \mathbf{p}^K$ , we evaluate the matrix spectral norm  $\|\mathcal{A} \times_3 \mathbf{p}^3 \cdots \times_K \mathbf{p}^K\|_2$  and take the maximum among them. Since the vertices of  $L_1$  balls are canonical basis vectors, this amounts to taking each matrix *slice* of  $\mathcal{A}$ , hence we call it the **slicing approach**. It yields the optimal solution for (9), and translates to an  $\alpha = \prod_{k=3}^K \sqrt{1/d_k}$  approximate solution for the TSN (6). The overall computational cost is  $O(\prod_{k=1}^K d_k)$ , which can be further reduced after embarrassing parallelization. This is much faster than existing matricization approaches [11–16], which cost  $O(\prod_{k=1}^K d_k \sum_{k=1}^K d_k)$  due to multiple *full* matrix SVDs.

A number of approximations of TSN is available in the literature [e.g., 26], all leading to the same  $\alpha$  guarantee here. We can easily combine them by picking the best one. Although the approximation ratio depends on

the problem size, it is only a worst case bound which is likely not tight. Perhaps more surprisingly, this is (up to logarithmic factors) the best result that we are aware of, despite being derived from such a straightforward polytopal approximation (more delicate tradeoff between computational cost and approximation quality is available in Appendix B). Significantly new ideas seem necessary for further improvement. In practice it is much more effective to perform block coordinate ascent (BCA) over (6) through all  $\mathbf{u}_k$ , initialized with the solution of (9). This is the strategy we use in our experiment. However, we also observed that randomly initializing BCA did not yield a result as good—a provably good initialization such as using (9) turns out to be crucial.

**Computational efficiency.** A key advantage of this scheme is that the computational efficiency can be further improved by utilizing the sparsity in  $\mathcal{A}$ . Suppose  $\mathcal{A}$  has  $z$  nonzeros entries. Then the optimization in (9) can be solved in  $\tilde{O}(z)$  time (independent of dimension), because the cost of computing the spectral norm of a matrix is linear in the number of nonzero entries. The subsequent BCA can also be performed in  $\tilde{O}(z)$  time per iteration, and [27] proposed a highly optimized algorithm.

## 4.2 Conditional Gradient with Approximate Spectral Norm

Equipped with an effective approximation of TSN, the challenge remains to design an *efficient* optimization algorithm which leverages this approximation guarantee and finds an  $\alpha$ -approximate solution to the original TTN regularized problem (8). We discover that this conversion can be accomplished by the recently developed generalized conditional gradient [GCG, 17, 18], which extends the traditional conditional gradient algorithm [28] to the (unconstrained) regularized problem. Technically, we further assume  $\ell$  is smooth (i.e. its gradient is Lipschitz continuous),

The basic GCG algorithm in [17] successively linearizes the loss  $\ell$  at the current iterate  $\mathcal{W}_t$ , finds an update direction in the unit ball of TTN:

$$\text{Oracle: } \mathcal{Z}_t \in \operatorname{argmin}_{\mathcal{Z}: \|\mathcal{Z}\|_{\text{tr}} \leq 1} \langle \mathcal{Z}, \nabla \ell(\mathcal{W}_t) \rangle, \quad (10)$$

and update by  $\mathcal{W}_{t+1} = (1 - \eta_t)\mathcal{W}_t + \eta_t \beta_t \mathcal{Z}_t$ , with some step size  $\eta_t \in [0, 1]$  and scaling factor

$$\beta_t = \operatorname{argmin}_{\beta \geq 0} \ell((1 - \eta_t)\mathcal{W}_t + \eta_t \beta \mathcal{Z}_t) + \lambda \cdot \eta_t \beta. \quad (11)$$

It is crucial to observe that the oracle problem (10) is exactly the same as the optimization involved in the definition of TSN (6). Therefore the above approximation of TSN readily provides an  $\alpha$ -approximate solution to the oracle problem here. So the final question

to resolve is whether or how GCG is “robust” against such approximate oracles. Theorem 1 provides an encouraging answer.

**Theorem 1.** *Let  $\ell \geq 0$  be convex, smooth, and have bounded sublevel sets. Denote  $f(\mathcal{W}) = \ell(\mathcal{W}) + \lambda \cdot \|\mathcal{W}\|_{\text{tr}}$ . Suppose in each iteration  $t$ , we find  $\mathcal{Z}_t$  that solves the oracle (10)  $\alpha$ -approximately, Then for all  $\mathcal{W}$  and for all  $t \geq 1$ , running GCG with  $\eta_t = \frac{2}{t+2}$  leads to  $f(\mathcal{W}_t) - \frac{f(\mathcal{W})}{\alpha} \leq \frac{2C}{t+3}$ , where  $C$  is some constant that does not depend on  $t$  or  $\alpha$ .*

The proof is relegated to Appendix C. Theorem 1 ensures that the GCG procedure is (asymptotically)  $\alpha$ -approximate when equipped with the  $\alpha$ -approximate oracle in (10). For instance, if  $\alpha = 1/2$ , then GCG is at most twice worse than the (possibly intractable) optimum. We note that [29] considered a different multiplicative approximation, requiring in each step an approximate solution of  $\min_{\|\mathcal{Z}\| \leq 1} \langle \mathcal{Z}, \nabla \ell(\mathcal{W}_t) \rangle - \langle \mathcal{W}_t, \nabla \ell(\mathcal{W}_t) \rangle$ . Unfortunately, this is usually hard to achieve due to the second changing term. On the other hand, [29] were still able to prove *exact* optimality while Theorem 1 here yields only an approximate guarantee. We noted in passing that a similar strategy appeared independently in [30] for hard matrix factorizations. In Appendix C, we further show that Theorem 1 can be generalized by replacing TTN with any convex positive homogeneous function.

Note this conversion is enabled by the convexity of  $f$ , which is facilitated by the use of TTN despite its own intractability. The polytopal approximation technique used for TSN, however, cannot be directly applied to the original problem (8), or in other numerical algorithms (e.g., proximal gradient). It is enabled by the linearization step (10) in GCG, where constraints in (6) decouple over modes and Euclidean balls can be approximated analytically and efficiently. Happily, any improvement on computing TSN immediately translates to the same amount for the tensor learning problem (8).

## 4.3 Local acceleration

A very efficient acceleration strategy was proposed in [17] for matrices. Pleasantly, we can extend their trick to the new tensor setting, starting with a variational form for the tensor trace norm:

**Theorem 2** (Variational formula). *Fix  $\mathcal{A} \in \mathbb{R}^{d_1 \times \dots \times d_K}$  and let  $t \geq \prod_{k=1}^K d_k$ . Then,*

$$\begin{aligned} \|\mathcal{A}\|_{\text{tr}} &= \min \left\{ \sum_{i=1}^t \|\mathbf{u}_i\|_2 \cdots \|\mathbf{z}_i\|_2 \right\} & (12) \\ &= \min \left\{ \frac{1}{K} \sum_{i=1}^t \|\mathbf{u}_i\|_2^K + \cdots + \|\mathbf{z}_i\|_2^K \right\}, & (13) \end{aligned}$$

where the minimum is taken w.r.t. all factorizations  $\mathcal{A} = \sum_{i=1}^t \mathbf{u}_i \otimes \cdots \otimes \mathbf{z}_i$ ,  $\mathbf{u}_i \in \mathbb{R}^{d_1}, \dots, \mathbf{z}_i \in \mathbb{R}^{d_K}$ .

---

**Algorithm 1:**  $\alpha$ -Approximate GCG for solving low-rank tensor learning (8) with local search
 

---

```

1 Set  $\mathcal{W}_0 = \mathbf{0}$ ,  $U_0 = \dots = Z_0 = [\ ]$ ,  $s_0 = 0$ .
2 for  $t = 1, 2, \dots$  do
3   find  $(\mathbf{u}_t, \dots, \mathbf{z}_t)$  that yields  $\alpha$ -approximate spectral norm of  $-\nabla\ell(\mathcal{W}_{t-1})$ ; // §4.1
4    $(a_t, b_t) \leftarrow \operatorname{argmin}_{a \geq 0, b \geq 0} \ell(a \cdot \mathcal{W}_{t-1} + b \cdot \mathbf{u}_t \otimes \dots \otimes \mathbf{z}_t) + \lambda(a \cdot s_t + b)$ ; // §4.2
5    $U_{\text{init}} \leftarrow (\sqrt[\kappa]{a_t} U_{t-1}, \sqrt[\kappa]{b_t} \mathbf{u}_t), \dots, Z_{\text{init}} \leftarrow (\sqrt[\kappa]{a_t} Z_{t-1}, \sqrt[\kappa]{b_t} \mathbf{z}_t)$ ;
6   find  $U_t, \dots, Z_t$  such that  $\mathfrak{F}_t(U_t, \dots, Z_t) \leq \mathfrak{F}_t(U_{\text{init}}, \dots, Z_{\text{init}})$ ; // §4.3
7    $\mathcal{W}_t \leftarrow \mathcal{I}_t \times_1 U_t \cdots \times_K Z_t$ ; // Implicitly maintained in implementation
8    $s_t \leftarrow \frac{1}{K} \sum_{i=1}^t \|(U_t)_{:i}\|_2^K + \dots + \|(Z_t)_{:i}\|_2^K$ ; // Upper bound on TTN, Thm 2
    
```

---

The first equality is well-known thanks to Grothendieck’s work on the tensor product of Banach spaces. Our proof in Appendix C gives an explicit bound on  $t$ , the number of rank-1 factors. In Appendix D we argue that the above variational forms should be preferred than some existing ones.

We now make use of Theorem 2 to further accelerate GCG. The entire procedure for solving the tensor learning problem (8) is summarized in Algorithm 1. After the  $(t-1)$ -th iteration of GCG (line 7), we get an explicit low-rank representation of the current iterate  $\mathcal{W}_{t-1} = \sum_{i=1}^{t-1} \mathbf{u}_i \otimes \dots \otimes \mathbf{z}_i$ . Instead of using the fixed step size  $\eta_t = \frac{2}{t+2}$  to combine  $\mathcal{W}_{t-1}$  and the newly found atom  $\mathbf{u}_t \otimes \dots \otimes \mathbf{z}_t$  (line 3) as in Theorem 1, we optimize it along with the scaling factor  $\beta$  (line 4, after change of variable to  $a, b$ ). To improve convergence, we locally minimize the surrogate function (line 6):

$$\mathfrak{F}_t(\tilde{U}, \dots, \tilde{Z}) := \ell(\sum_{i=1}^t \tilde{U}_{:i} \otimes \dots \otimes \tilde{Z}_{:i}) + \frac{\lambda}{K} \sum_{i=1}^t \|\tilde{U}_{:i}\|_2^K + \dots + \|\tilde{Z}_{:i}\|_2^K, \quad (14)$$

simply replacing the intractable TTN with the variational form (12). We expect the surrogate problem (with fixed  $t$ ) to be reasonably close to the original problem (8), hence a local (smooth) unconstrained minimization of (14) should improve the iterate  $\mathcal{W}_{t-1}$  of GCG. Adopting proper initialization to ensure monotonic progress (line 5), the convergence guarantee in Theorem 1 is retained.

Although local optimization does not improve convergence rate in the worst case, empirically it much accelerates convergence by allowing the factors found in past iterations to be re-optimized in the context of new ones. This provides much more freedom than merely optimizing their weights as in the vanilla GCG. Faster convergence also results in simpler models which improves generalization performance in general. Interestingly, GCG with local search bears much resemblance to ALS, but with the key difference being that the number of factors is not fixed, and new factors are incorporated in a greedy fashion. This potentially alleviates the issue of local optimum suffered by ALS.

**Computational efficiency.** Note in Algorithm 1,  $\mathcal{W}_t$  is always represented via CP decomposition  $\mathcal{W}_t = \sum_{i=1}^t \mathbf{u}_i \otimes \dots \otimes \mathbf{z}_i$ , which is cheap to maintain and underpins the seamless integration of GCG and local search. The most costly step in Algorithm 1 is computing the gradients  $\frac{\partial \ell}{\partial \mathbf{u}_i}, \dots, \frac{\partial \ell}{\partial \mathbf{z}_i}$  given  $\nabla \ell(\mathcal{W}_{t-1})$ , for which a naive approach takes  $O(K \prod_{k=1}^K d_k)$  time. We show in Appendix E that this cost can be reduced to  $O(\prod_{k=1}^K d_k)$  via dynamic programming.

A key merit of GCG lies in its ability to effectively utilize both the sparsity in tensor completion, and the low rank of the optimal solution. Assuming the tensor has  $z$  nonzeros, the approximate oracle costs  $O(z)$  time, and at iteration  $t$  the gradients in local search cost  $O(zt)$  [27], on par with ALS.  $t$  is typically low since the optimal solution has low rank, allowing GCG to converge fast. As a result, its overall cost is much cheaper than matricization based methods, which require  $O(z \sum_k d_k)$  computation for singular value thresholding on a sparse matrix that is *not* necessarily low-rank in intermediate steps. Finally, the space cost is only  $O(z + t \sum_k d_k)$ , much less than the  $O(\prod_k d_k)$  of matricization approaches.

## 5 Sample Complexity Comparisons

Arguably, for any tractable approximation of TSN, we can directly use its dual in the optimization objective (8), in lieu of TTN. This raises the question of why not directly analyze their sample complexity, and why do we analyze its deviation from TTN. This is because using TTN as a gold standard allows us to infer *suggestively* that a solution based on tighter approximation of TSN will lead to as good or better sample complexity. For example, once we establish the sample complexity for the slicing algorithm (below), then intuitively it should not deteriorate when conjoined with BCA, which improves the approximation of TSN. We leave a rigorous analysis for future work.

To give some theoretical underpinnings, let us consider

concretely the tensor completion problem:

$$\min_{\mathcal{W}} \sum_{i=1}^q \|\mathcal{L}_i(\mathcal{W})\|_{(i)} \quad \text{s.t.} \quad \langle \mathcal{A}_i, \mathcal{W} \rangle = b_i, \quad (15)$$

where each  $\mathcal{A}_i$  is a Gaussian random tensor, e.g., the entries of  $\mathcal{A}_i$  are i.i.d. standard normal, and  $b_i = \langle \mathcal{A}_i, \mathcal{W}_0 \rangle$  for some unknown low-rank tensor  $\mathcal{W}_0$ . For simplicity, let  $d_1 = \dots = d_K = d$ . The linear maps  $\mathcal{L}_i$  are from  $\mathbb{R}^d \times \dots \times \mathbb{R}^d$  to  $\mathbb{R}^{m_i}$ , and the norm  $\|\cdot\|_{(i)}$  is defined on the corresponding range space  $\mathbb{R}^{m_i}$ . In particular, we will compare the following two approaches:

- **Matricization:**  $q = K$ , and for  $i = 1, \dots, q$ ,  $\mathcal{L}_i : \mathbb{R}^d \times \dots \times \mathbb{R}^d \rightarrow \mathbb{R}^d \times \mathbb{R}^{d^{K-1}}$ ,  $\mathcal{W} \mapsto \mathcal{W}_{(i)}$ , and  $\|\cdot\|_{(i)}$  is the matrix trace norm on  $\mathbb{R}^d \times \mathbb{R}^{d^{K-1}}$ .
- **Slicing:**  $q = d^{K-2}$ , and for  $i = 1, \dots, q$ ,  $\mathcal{L}_i : \mathbb{R}^d \times \dots \times \mathbb{R}^d \rightarrow \mathbb{R}^d \times \mathbb{R}^d$ ,  $\mathcal{W} \mapsto \mathcal{W}_{:,i_3,\dots,i_K}$ , and  $\|\cdot\|_{(i)}$  is the matrix trace norm on  $\mathbb{R}^d \times \mathbb{R}^d$ , where  $i = \sum_{j=3}^K d^{j-3} i_j$  is the  $d$ -ary expansion.

The matricization case has been widely studied in [e.g. 11–15, 31], while the slicing case does not appear to be widely appreciated. We used slicing in Section 4.1 to approximate the TSN.

We are interested in deciding the least number of observations  $m$  that still allows recovering  $\mathcal{W}_0$  (with high probability) by solving (15). Obviously, the recovery is successful iff  $\mathcal{W}_0$  is the unique minimizer of (15). Let us recall the following result:

**Theorem 3** ([31, 32]). *Suppose  $\mathcal{W}_0 \neq \bigcap_{i \leq q} \text{null}(\mathcal{L}_i)$ . For each  $i$ , define  $L_i = \sup_{0 \neq \mathbf{x} \in \mathbb{R}^{m_i}} \frac{\|\mathbf{x}\|_{(i)}}{\|\mathbf{x}\|_F}$ . Set  $\kappa_i = \frac{d^K \|\mathcal{L}_i(\mathcal{W}_0)\|_{(i)}^2}{L_i^2 \|\mathcal{L}_i\|^2 \|\mathcal{W}_0\|_F^2}$ . Then for  $m \leq O(\min_i \kappa_i)$ , with high probability  $\mathcal{W}_0$  is not the unique minimizer of (15), while for  $m \geq O(\max_i \kappa_i)$ , with high probability  $\mathcal{W}_0$  is the unique minimizer.*

We apply this theorem to the matricization and the slicing approach above. Note first that in both cases  $L_i = \sqrt{d}$  and the operator norm  $\|\mathcal{L}_i\| = 1$ . Thus, for successful recovery we need to compare:

$$\frac{d^{K-1}}{\|\mathcal{W}_0\|_F^2} \max_i \|(\mathcal{W}_0)_{(i)}\|_{\text{tr}}^2 \quad \text{vs.} \quad \frac{d^{K-1}}{\|\mathcal{W}_0\|_F^2} \max_i \|(\mathcal{W}_0)_{:,i_3,\dots,i_K}\|_{\text{tr}}^2.$$

Clearly, both are on the same order  $O(d^{K-1})$  but we always have  $\max_i \|(\mathcal{W}_0)_{(i)}\|_{\text{tr}}^2 \geq \max_i \|(\mathcal{W}_0)_{:,i_3,\dots,i_K}\|_{\text{tr}}^2$ , i.e., the slicing approach (latter) always has a smaller problem dependent constant. This can make a big difference: We could obtain the correct order  $O(d^{K-1})$  by using say even the  $\ell_1$  norm with  $\mathcal{L}_i : \mathcal{W} \mapsto \mathcal{W}_{:,i_2,\dots,i_K}$ . Only the problem dependent term inside the max function reflects how different norms match the problem structure (e.g. low-rank). We note that by using a more balanced unfolding, [31] was able to improve the

order  $O(d^{K-1})$ , at the expense of potentially increasing the problem dependent term. On the other hand, the slicing norm is particularly computational friendly: it can be embarrassingly parallelized. Moreover, when  $\mathcal{A}_i$  are sampling operators, the problem becomes completely separable and we need only complete each slice independently.

We wish to point out a second difference between the matricization approach and the proposed  $\alpha$ -approximate GCG with the slicing norm: The latter not only recovers a low-rank tensor, it also finds a CP decomposition that certifies the low-rank assumption. In contrast, the matricization approach can only recover the tensor, but it never explicitly finds a usable CP decomposition.

A large body of works deal directly with the tensor rank function, e.g. [1, 33, 34]. Nonconvex optimization tools are employed to get a local minima, whose formal guarantee, to our best knowledge, is largely unknown. Under the restrictive orthogonal CP assumption (which may not exist in general), [4] justified an extended form of power iteration while [2] proved strong guarantees by restricting the rank. It may be possible to strengthen our results by combining ideas from them.

## 6 Experiments

We study the empirical performance of Algorithm 1 on two low-rank learning problems: tensor completion and multitask learning. Local optimization (14) in GCG was solved by LBFGS. ALS was tested in all cases, with gradients computed by the state-of-the-art algorithm for efficiency [27].

### 6.1 Low-Rank Tensor Completion

Two state-of-the-art algorithms were used for comparison: HaLRTC [14] which uses the weighted matrix trace norm; and a tighter relaxation by [35], referred to as RpLRTC. [14, Eq 11] also considered a non-convex Parafac regularizer that is similar to the variational form of tensor trace norm. Since its performance has been shown inferior to HaLRTC, we will not include it in the comparison. Also note the analysis of [4] does not consider our noisy setting. For all methods, we selected the value of  $\lambda$  via a validation set, which always consisted of 10% of the whole dataset. The rank in ALS was set to the true value on synthetic data, and was selected via a validation set on real data.

**Synthetic data.** We first generated a 3-order tensor  $\mathcal{Z}^0 \in \mathbb{R}^{d \times d \times d}$  using the CP decomposition (1), where  $\mathbf{u}_i, \mathbf{v}_i, \mathbf{z}_i$  were drawn i.i.d. from the standard normal

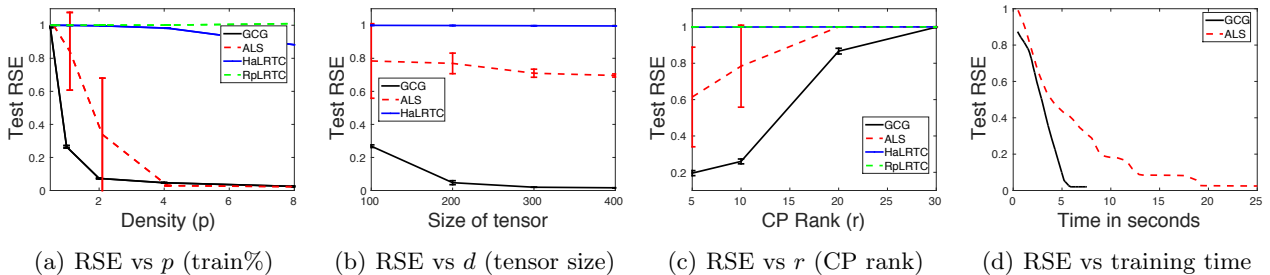


Figure 1: Test RSE on synthetic data generated by CP decomposition. RpLRTC is not included in (b) because it gets very slow when  $d \geq 200$ .

distribution. So  $\text{rank}(\mathcal{Z}^0) \leq r$  and we loosely call  $r$  the rank.  $\mathcal{Z}_{i,j,k} := (\mathcal{Z}_{i,j,k}^0 - \text{mean}(\mathcal{Z}^0)) / (d^{1.5} \text{std}(\mathcal{Z}^0)) + \xi_{i,j,k}$ , where  $\text{mean}$  and  $\text{std}$  stand for the mean and standard deviation, respectively, and  $\xi_{i,j,k}$  are i.i.d. Gaussians with zero mean and variance  $0.1^2$ . We randomly picked  $p\%$  elements of  $\mathcal{Z}$  as observations for training, and the reconstructed tensor  $\mathcal{W}$  was compared with  $\mathcal{Z}^0$  on the rest  $(90 - p)\%$  entries. The whole process was repeated 10 times, and we report the average of  $\text{RSE} := \|\mathcal{W}_{\text{te}} - \mathcal{Z}_{\text{te}}^0\|_F / \|\mathcal{Z}_{\text{te}}^0\|_F$  [14].

Figures 1(a) to 1(c) show the test RSE as a function of density  $p$  of the training data, the size  $d$ , and the CP rank  $r$  respectively. In all plots, parameters not varying in the  $x$ -axis were set to the default values:  $p = 1$ ,  $r = 10$ , and  $d = 100$ . Clearly, in all cases GCG yields significantly more accurate completions than ALS, HaLRTC, and RpLRTC. Fixing  $d$  and  $r$ , GCG clearly outperforms when the tensor is sparse ( $p = 1$  or 2). In real datasets such as Yelp and Nell below, less than 1% entries are observed, and therefore it is indeed an interesting regime. ALS exhibits higher variance suggesting susceptibility to local minima. With fixed  $p$  and  $r$ , increasing the size of tensor allows both ALS and GCG to reduce test RSE. But it benefits GCG much more thanks to its convexity. Finally, increasing the CP rank makes the problem challenging for all methods. This is reasonable because the sample size is fixed. Overall, the results confirm higher faithfulness of GCG in low-rank tensor completion than nonconvex methods and matricization based regularizers when the tensor is sparse.

Computational time is another key aspect of comparison. To make the comparison fair, we set  $d = 100$ ,  $p = 8$ ,  $r = 10$  so that ALS and GCG achieve similar RSE. HaLRTC and RpLRTC are not included because they scale much more poorly. Figure 1(d) shows that GCG is much more efficient than ALS in reducing the test RSE.

**Image inpainting.** Images are naturally represented as a “width (259)”  $\times$  “height (247)”  $\times$  “RGB

(3)” tensor. We used an image of facade [14] and randomly sampled  $p\%$  pixels from the original color image for training. All pixel values were rescaled to  $[0, 1]$ . In addition to the competing methods used above, we also compared with the latent Schatten norm [15] as the dataset does not necessarily have low CP rank any more. Figure 2 reports the RSE on the remaining  $(90 - p)\%$  pixels with varied  $p$ . Clearly, GCG is superior in propagating global structure from a small number of observed pixels. Again ALS exhibits high variance, probably due to local minima.

**Multitask recommendation.** We next study tensor completion for recommendation systems, and the focus is on large-scale real data. The Yelp dataset is from the Yelp Data Challenge [36] which contains (user, items, word) triples for business reviews. It is sized  $46K \times 12K \times 85K$  with 9.9 million nonzeros. The Nell dataset is from [37] which consists of (noun phrase 1, context, noun phrase 2) triples. It is sized  $12K \times 9K \times 29K$  with 77 million nonzeros, and we randomly subsampled 10% for experiment. None of the matricization based methods can scale to such large datasets, and therefore we only compared GCG with ALS. We varied the number of training (observed) entries, used another 10% entries for validation, and tested on the rest entries. The random partition was repeated for 10 times, and the mean test RSE is presented in Figure 3 and 4. For a range of training set size, GCG yields significantly lower error than ALS.

## 6.2 Multitask Learning

**Restaurant recommendation [13].** The tasks are to predict the rating that a restaurant would receive from each of the  $I_2 = 138$  customers, in  $I_3 = 3$  aspects (food, service, overall). So there are  $I_2 \times I_3$  tasks, with each task employing a  $I_1 = 45$  dimensional weight vector. The  $I_1 \times I_2 \times I_3$  weight tensor is assumed to have low rank. Each of the 3483 data points is a restaurant represented by an  $I_1$  dimensional feature vector (e.g. location and cuisine type), as well as its

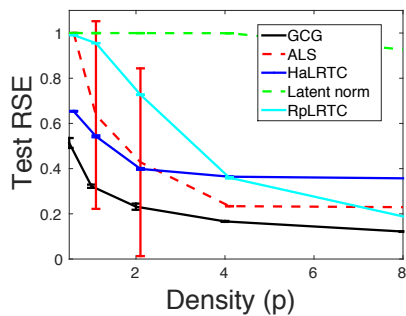


Figure 2: Image inpainting

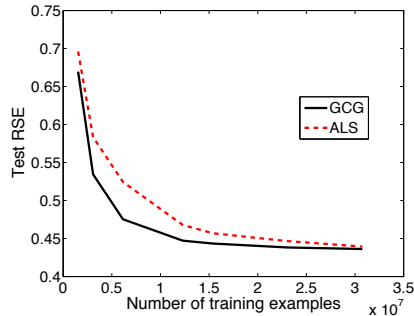


Figure 3: Test RSE on Yelp

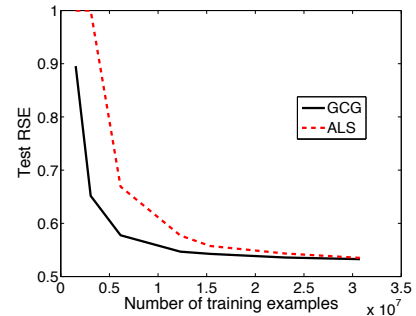


Figure 4: Test RSE on Nell

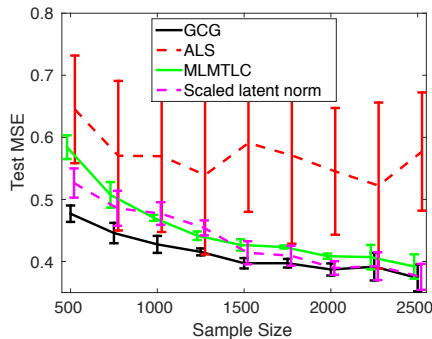


Figure 5: Restaurant rating prediction

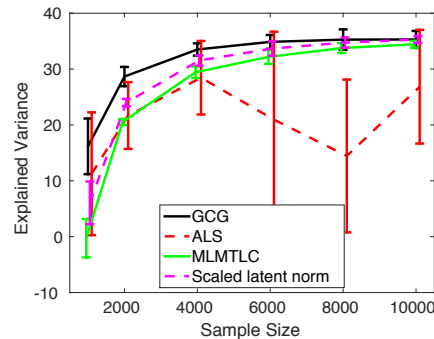


Figure 6: School grade prediction

$I_2 \times I_3$  number of ratings. Predicting these ratings constitutes  $I_2 \times I_3$  regression problems, each with an  $\ell_2$  loss.

GCG was compared with ALS, a convex multilinear multitask learning model (MLMTL-C) [13], and scaled latent trace norm proposed recently in [38]. They have been shown to predict more accurately than a number of other multitask algorithms (hence not included here). We randomly sampled  $m$  ratings (across all tasks) for training, 250 ratings for validation (used by GCG only), and the rest for testing. All competing algorithms selected models by optimizing the test mean squared error (MSE), an extra advantage that was not available to GCG. Figure 5 plots the test MSE as a function of  $m$ , averaged over 20 random repetitions. Clearly, GCG yields significantly lower test error over a range of training set sizes, with the edge diminishing as the sample size grows. It is also more efficient. For example, it took 0.5 seconds for GCG to train on 1500 examples, while MLMTL-C and scaled latent norm spent about 9 seconds and 12 seconds respectively.

**School grade prediction.** Following [38], we used multitask learning to model the Inner London Education Authority dataset. It records the score of 15,362 students at 139 schools over 3 years, with each student represented by 24 attributes. The weight tensor is  $24 \times 139 \times 3$ , and following [38] we measured the

percentage of explained variance:  $100 \cdot (1 - (\text{test MSE}) / (\text{variance of score}))$ . Model selection in GCG was based on a validation set of 1000 samples, while all competing algorithms were allowed to select models by directly optimizing the test measure. As shown in Figure 6, GCG again achieves the highest performance, and is particularly advantageous when the sample size is small (2000 to 6000). The results of MLMTL-C and scaled latent norm recover those in [38].

## 7 Conclusion

Learning a low-rank tensor is an important task in many real-world applications. In this work we consider the genuine tensor trace norm regularization, as a proxy of the nonconvex tensor rank function. Building on a simple polytopal approximation of the dual spectral norm, we demonstrate how to solve the tensor trace norm problem using the recent GCG algorithm. We establish the  $O(1/t)$  rate of convergence of GCG with such an approximate dual spectral norm. Further accelerations using an unconstrained smooth surrogate is discussed. Compared to matricization approaches, our algorithm achieves the same order of sample complexity but with potentially a smaller problem dependent constant. Computationally, our algorithm benefits from both the sparse and low-rank structure hence scales much better. Experiments on a variety of real datasets confirm the effectiveness of our algorithm.



## References

- [1] T. G. Kolda and B. W. Bader. Tensor decompositions and applications. *SIAM Review*, 51:455–500, 2009.
- [2] A. Anandkumar, R. Ge, and M. Janzamina. Guaranteed non-orthogonal tensor decomposition via alternating rank-1 updates, 2014. ArXiv:1402.5180v2.
- [3] M. Mørup and L. K. Hansen. Automatic relevance determination for multi-way models. *Journal of Chemometrics*, 23(7-8):352–363, 2009.
- [4] P. Jain and S. Oh. Provable tensor factorization with missing data. In *NIPS*. 2014.
- [5] V. Kuleshov, A. Chaganty, and P. Liang. Tensor factorization via matrix factorization. In *AISTATS*. 2015.
- [6] E. J. Candès and B. Recht. Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9:717–772, 2009.
- [7] V. Chandrasekaran, B. Recht, P. A. Parrilo, and A. S. Willsky. The convex geometry of linear inverse problems. *Foundations of Computational Mathematics*, 12(6):805–849, 2012.
- [8] M. Yuan and C.-H. Zhang. On tensor completion via nuclear norm minimization. *Foundations of Computational Mathematics*, to appear.
- [9] V. D. Silva and L.-H. Lim. Tensor rank and the ill-posedness of the best low-rank approximation problem. *SIAM Journal on Matrix Analysis and Applications*, 30(3):1084–1127, 2008.
- [10] C. J. Hillar and L.-H. Lim. Most tensor problems are NP-hard. *Journal of ACM*, 60(6):45:1–45:39, 2013.
- [11] M. Signoretto, L. D. Lathauwer, and J. A. K. Suykens. Nuclear norms for tensors and their use for convex multilinear estimation. *Tech. Rep. 10-186*, ESAT-SISTA, K. U. Leuven, 2010.
- [12] S. Gandy, B. Recht, and I. Yamada. Tensor completion and low-n-rank tensor recovery via convex optimization. *Inverse Problems*, 27:1–19, 2011.
- [13] B. Romera-Paredes, M. Aung, N. Bianchi-Berthouze, and M. Pontil. Multilinear multitask learning. In *ICML*. 2013.
- [14] J. Liu, P. Musialski, P. Wonka, and J. Ye. Tensor completion for estimating missing values in visual data. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(1):208–220, 2013.
- [15] R. Tomioka and T. Suzuki. Convex tensor decomposition via structured Schatten norm regularization. In *NIPS*. 2013.
- [16] D. Goldfarb and Z. Qin. Robust low-rank tensor recovery: Models and algorithms. *SIAM Journal on Matrix Analysis and Applications*, 35:225–253, 2014.
- [17] X. Zhang, Y. Yu, and D. Schuurmans. Accelerated training for matrix-norm regularization: A boosting approach. In *NIPS*. 2012.
- [18] Z. Harchaoui, A. Juditsky, and A. Nemirovski. Conditional gradient algorithms for norm-regularized smooth convex optimization. *Mathematical Programming*, 152:75–112, 2015.
- [19] A. Stegeman and P. Comon. Subtracting a best rank-1 approximation may increase tensor rank. *Linear Algebra and its Applications*, 433:1276–1300, 2010.
- [20] M. Fazel, H. Hindi, and S. P. Boyd. A rank minimization heuristic with application to minimum order system approximation. In *American Control Conference*, pp. 4734–4739. 2001.
- [21] L.-H. Lim and P. Comon. Blind multilinear identification. *IEEE Transactions on Information Theory*, 60(2):1260–1280, 2014.
- [22] H. Derksen. On the nuclear norm and the singular value decomposition of tensors. *Foundations of Computational Mathematics*, to appear.
- [23] N. Rao, P. Shah, and S. Wright. Forward - backward greedy algorithms for atomic norm regularization, 2014. ArXiv:1404.5692.
- [24] M. Kochol. Constructive approximation of a ball by polytopes. *Math Slovaca*, 44(1):99–105, 1994.
- [25] I. Bárány and Z. Füredi. Approximation of the sphere by polytopes having few vertices. *Proceedings of the American Mathematical Society*, 102(3):651–659, 1988.
- [26] Z. Li, S. He, and S. Zhang. *Approximation Methods for Polynomial Optimization: Models, Algorithms, and Applications*. Springer, 2012.
- [27] J. H. Choi and S. Vishwanathan. Dfacto: Distributed factorization of tensors. In *NIPS*. 2014.
- [28] M. Frank and P. Wolfe. An algorithm for quadratic programming. *Naval Research Logistics Quarterly*, 3(1-2):95–110, 1956.
- [29] S. Lacoste-Julien, M. Jaggi, M. Schmidt, and P. Pletscher. Block-Coordinate Frank-Wolfe optimization for structured svms. In *ICML*. 2013.
- [30] F. Bach. Convex relaxations of structured matrix factorizations, 2013. HAL:00861118.
- [31] C. Mu, B. Huang, J. Wright, and D. Goldfarb. Square deal: Lower bounds and improved relaxations for tensor recovery. In *ICML*. 2014.

- [32] D. Amelunxen, M. Lotz, M. B. McCoy, and J. A. Tropp. Living on the edge: phase transitions in convex programs with random data. *Information and Inference*, 3(3):224–294, 2014.
- [33] J. Nie and L. Wang. Semidefinite relaxations for best rank-1 tensor approximations. *SIAM Journal on Matrix Analysis and Applications*, 35(3):1155–1179, 2014.
- [34] B. Jiang, S. Ma, and S. Zhang. Tensor principal component analysis via convex optimization. *Mathematical Programming*, 150:423–457, 2015.
- [35] B. Romera-Paredes and M. Pontil. A new convex relaxation for tensor completion. *In NIPS*. 2013.
- [36] [https://www.yelp.com/dataset\\_challenge/dataset](https://www.yelp.com/dataset_challenge/dataset).
- [37] U. Kang, E. Papalexakis, A. Harpale, and C. Faloutsos. Gigatensor: Scaling tensor analysis up by 100 times - algorithms and discoveries. *In KDD*. 2012.
- [38] K. Wimalawarne, M. Sugiyama, and R. Tomioka. Multitask learning meets tensor factorization: task imputation via convex optimization. *In NIPS*. 2014.
- [39] S. Becker, J. Bobin, and E. J. Candès. NESTA: A fast and accurate first-order method for sparse recovery. *SIAM Journal on Imaging Sciences*, 4(1):1–39, 2009.
- [40] R. Bro. PARAFAC: Tutorial and applications. *Chemometrics and Intelligent Laboratory Systems*, 38:147–171, 1997.
- [41] S. Engelen, S. Frosch, and B. Jorgensen. A fully robust parafac method for analyzing fluorescence data. *Journal of Chemometrics*, 23(3-4):124–131, 2009.

## A Matricization

A tensor  $\mathcal{A}$  can be unfolded into a matrix  $A$  in various ways. We focus here on 2-way unfoldings that specify a proper partition of  $\{1, \dots, K\} := [K] = \{m_1, \dots, m_p\} \cup \{n_1, \dots, n_q\}$ , an integer-valued row index function  $\mathbf{r}$ , and a column index function  $\mathbf{c}$ , such that  $\mathcal{A}_{i_1, \dots, i_K} \mapsto A_{\mathbf{r}(i_{m_1}, \dots, i_{m_p}), \mathbf{c}(i_{n_1}, \dots, i_{n_q})}$ . This is simply a rearrangement of the entries in  $\mathcal{A}$  hence preserves its Frobenius norm. We use the notation  $\mathcal{A}_{(\mathbf{r}, \mathbf{c})}$  for the 2-way unfolding under the index functions  $\mathbf{r}$  and  $\mathbf{c}$ .

The above tensor unfolding interacts conveniently with the mode- $k$  multiplication, once we define a suitable matrix Kronecker product. Fix two row index functions  $\mathbf{r}$  and  $\hat{\mathbf{r}}$ . The Kronecker product of  $p$  matrices  $U \in \mathbb{R}^{\hat{d}_{m_1} \times d_{m_1}}, \dots, W \in \mathbb{R}^{\hat{d}_{m_p} \times d_{m_p}}$  is a matrix that has size  $\prod_{k=1}^p \hat{d}_{m_k} \times \prod_{k=1}^p d_{m_k}$  and satisfies

$$(U \boxtimes \dots \boxtimes W)_{\hat{i}, i} = U_{\hat{i}_{m_1}, i_{m_1}} \dots W_{\hat{i}_{m_p}, i_{m_p}}, \quad (16)$$

where  $\hat{i} = \hat{\mathbf{r}}(\hat{i}_{m_1}, \dots, \hat{i}_{m_p})$  and  $i = \mathbf{r}(i_{m_1}, \dots, i_{m_p})$ .

Similar definitions can be made using two column index functions  $\mathbf{c}$  and  $\hat{\mathbf{c}}$ . It is just algebra to verify that

$$(\mathcal{A} \times_1 U_1 \dots \times_K U_K)_{(\hat{\mathbf{r}}, \hat{\mathbf{c}})} = (U_{m_1} \boxtimes \dots \boxtimes U_{m_p}) \mathcal{A}_{(\mathbf{r}, \mathbf{c})} (U_{n_1} \boxtimes \dots \boxtimes U_{n_q})^\top, \quad (17)$$

where the first and second group of Kronecker product use  $(\hat{\mathbf{r}}, \mathbf{r})$  and  $(\hat{\mathbf{c}}, \mathbf{c})$  respectively.

**Example 1** (Mode- $k$  unfolding  $\mathcal{A}_{(k)}$ ). *To illustrate the above definition, let us consider the partition  $[K] = \{k\} \cup \{1, \dots, k-1, k+1, \dots, K\}$  and the index functions*

$$\mathbf{c}(i_1, \dots, i_{k-1}, i_{k+1}, \dots, i_K) = 1 + \sum_{j \neq k} (i_j - 1) \prod_{m > j, m \neq k} d_m,$$

and  $\mathbf{r}(i_k) = i_k$ . This is called the mode- $k$  unfolding, together with the notation  $\mathcal{A}_{(k)}$ . Here  $\boxtimes$  reduces to the usual matrix Kronecker product, and  $(\mathcal{A} \times_k U)_{(k)} = U \mathcal{A}_{(k)}$ . For matrices, we simply have  $A_{(1)} = A$ ,  $A_{(2)} = A^\top$ , and  $A \times_1 U \times_2 V = U A V^\top$ .

**Example 2** (Balanced mode- $k$  unfolding  $\mathcal{A}_{[k]}$ ). *The mode- $k$  unfolding above yields an extremely unbalanced matrix with size  $d_k \times \prod_{j \neq k} d_j$ . A more balanced unfolding is proposed in [31], consisting of the partition  $[K] = \{1, \dots, k\} \cup \{k+1, \dots, K\}$  and the index functions*

$$\begin{aligned} \mathbf{r}(i_1, \dots, i_k) &= 1 + \sum_{j=1}^k (i_j - 1) \prod_{k \geq m \geq j+1} d_m \\ \mathbf{c}(i_{k+1}, \dots, i_K) &= 1 + \sum_{j=k+1}^K (i_j - 1) \prod_{m \geq j+1} d_m, \end{aligned}$$

with  $\hat{\mathbf{r}}, \hat{\mathbf{c}}$  similarly defined. The resulting matrix, denoted as  $\mathcal{A}_{[k]}$ , has size  $\prod_{j=1}^k d_j \times \prod_{j=k+1}^K d_j$ . For  $k = \lfloor K/2 \rfloor$ , the unfolding is more like a square matrix, which can be beneficial in completion tasks [31].

## B Approximating Tensor Spectral Norm

Let  $\mathbb{B}_2^d$  be the Euclidean norm ball in a  $d$  dimensional space ( $d$  is a superscript in  $\mathbb{B}_2^d$ ). We can approximate  $\mathbb{B}_2^d$  with a *polytope*, based on a celebrated result from convex geometry [24]: For any  $d, n \geq 2$  we can find in polynomial time (at most)  $d^n$  points  $\{\mathbf{p}_i\}$ , such that their convex hull  $\mathbb{P}^d$  satisfies

$$\mathbb{P}^d \subseteq \mathbb{B}_2^d \subseteq \frac{1}{c} \sqrt{\frac{d}{n \log d}} \mathbb{P}^d, \quad (18)$$

where  $c$  is some universal constant. For  $n = 1$ , we can simply take the unit ball of the 1-norm, denoted as  $\mathbf{B}_1^d$ , and get a similar result

$$\mathbf{B}_1^d \subseteq \mathbf{B}_2^d \subseteq \sqrt{d}\mathbf{B}_1^d. \quad (19)$$

By counting the volume, it can be proved that the factors in (18) and (19) are the best possible respectively [25]. Specializing to the tensor spectral norm, we simply replace each Euclidean ball constraint with its polytopal approximation as suggested in (18) or (19), and evaluate the inner product at each of the vertices of the polytope. Since the matrix trace norm is tractable, we need only execute the polytopal approximation for the last  $K - 2$  balls, i.e., solving the approximation

$$\max_{\mathbf{u}_1 \in \mathbf{B}_2^{d_1}, \mathbf{u}_2 \in \mathbf{B}_2^{d_2}, \mathbf{u}_3 \in \mathbf{P}^{d_3}, \dots, \mathbf{u}_K \in \mathbf{P}^{d_K}} \langle \mathcal{A}, \mathbf{u}_1 \otimes \dots \otimes \mathbf{u}_K \rangle. \quad (20)$$

So for each vertex  $\mathbf{p}^3 \otimes \dots \otimes \mathbf{p}^K$  (with  $\prod_{k=3}^K (d_k)^n$  of them), we evaluate the matrix spectral norm  $\|\mathcal{A} \times_3 \mathbf{p}^3 \cdots \times_K \mathbf{p}^K\|_2$  and pick the maximum. This yields the optimal solution for (20), which immediately translates to an  $\alpha = O(\prod_{k=3}^K \sqrt{nd_k^{-1} \log d_k})$  approximate solution for the tensor spectral norm (6). If we set  $n = 1$  and use (19), then  $\alpha = \prod_{k=3}^K \sqrt{1/d_k}$  and only  $\prod_{k=3}^K d_k$  matrix spectral norms need to be checked. The overall computational cost is  $O(\prod_{k=1}^K d_k)$ . It is also easy to reduce each factor  $d_k$  to the smaller constant rank( $\mathcal{A}_{(k)}$ ), or simply rank( $\mathcal{A}$ ). For  $n \geq 2$  we get a  $\log d_k$  factor improvement in the approximation guarantee, at the expense of a more complicated and costly implementation.

## C Proofs omitted in Section 4

**Theorem 1.** *Let  $\ell \geq 0$  be convex, smooth, and have bounded sublevel sets. Denote  $f(\mathcal{W}) = \ell(\mathcal{W}) + \lambda \cdot \kappa(\mathcal{W})$ . Suppose in each iteration  $t$ , we find  $\mathcal{Z}_t$  that satisfies*

$$\kappa(\mathcal{Z}_t) \leq 1, \quad \langle \mathcal{Z}_t, \nabla \ell(\mathcal{W}_t) \rangle \geq \alpha \cdot \max_{\kappa(\mathcal{Z}) \leq 1} \langle \mathcal{Z}, \nabla \ell(\mathcal{W}_t) \rangle. \quad (21)$$

*Then for all  $\mathcal{W}$  and for all  $t \geq 1$ , running GCG with  $\eta_t = \frac{2}{t+2}$  leads to  $f(\mathcal{W}_t) - \frac{f(\mathcal{W})}{\alpha} \leq \frac{2C}{t+3}$ , where  $C$  is some constant that does not depend on  $t$  or  $\alpha$ .*

Recall that the function  $\ell$  is smooth if its gradient is Lipschitz continuous with respect to some norm  $\|\cdot\|$ , namely that for all  $\mathcal{W}$  and  $\mathcal{Z}$ ,

$$\ell(\mathcal{Z}) \leq \ell(\mathcal{W}) + \langle \mathcal{Z} - \mathcal{W}, \nabla \ell(\mathcal{W}) \rangle + \frac{L}{2} \|\mathcal{Z} - \mathcal{W}\|^2, \quad (22)$$

for some constant  $L := L_{\|\cdot\|} \geq 0$ . The least squares loss in fact satisfies (22) with equality and  $L = 2$ .

*Proof.* Let  $\mathcal{W}$  be arbitrary and  $s = \kappa(\mathcal{W})$ . Let  $\tilde{\mathcal{W}}_{t+1}$  be the output of GCG at iteration  $t + 1$  and  $\mathcal{W}_{t+1}$  be the

improved iterate after local search. The following chain of inequalities can be easily verified:

$$\begin{aligned}
 f(\mathcal{W}_{t+1}) &\leq f(\tilde{\mathcal{W}}_{t+1}) \\
 &= \ell(\tilde{\mathcal{W}}_{t+1}) + \lambda \cdot \kappa(\tilde{\mathcal{W}}_{t+1}) \\
 &= \ell((1 - \eta_t)\mathcal{W}_t + \eta_t\beta_t\mathcal{Z}_t) + \lambda \cdot \kappa((1 - \eta_t)\mathcal{W}_t + \eta_t\beta_t\mathcal{Z}_t) && \text{(definition of } \tilde{\mathcal{W}}_{t+1}\text{)} \\
 &\leq \ell((1 - \eta_t)\mathcal{W}_t + \eta_t\beta_t\mathcal{Z}_t) + \lambda(1 - \eta_t)\kappa(\mathcal{W}_t) + \lambda\eta_t\beta_t\kappa(\mathcal{Z}_t) && \text{(sublinearity of } \kappa\text{)} \\
 &\leq \ell((1 - \eta_t)\mathcal{W}_t + \eta_t\beta_t\mathcal{Z}_t) + \lambda(1 - \eta_t)\kappa(\mathcal{W}_t) + \lambda\eta_t\beta_t && \text{(definition of } \mathcal{Z}_t\text{)} \\
 &\leq \ell\left((1 - \eta_t)\mathcal{W}_t + \eta_t\frac{s}{\alpha}\mathcal{Z}_t\right) + \lambda(1 - \eta_t)\kappa(\mathcal{W}_t) + \lambda\eta_t\frac{s}{\alpha} && \text{(definition of } \beta_t\text{)} \\
 &\leq f(\mathcal{W}_t) + \eta_t \left\langle \frac{s}{\alpha}\mathcal{Z}_t - \mathcal{W}_t, \nabla\ell(\mathcal{W}_t) \right\rangle + \frac{L \left\| \frac{s}{\alpha}\mathcal{Z}_t - \mathcal{W}_t \right\|^2}{2} \eta_t^2 - \lambda\eta_t\kappa(\mathcal{W}_t) + \lambda\eta_t\frac{s}{\alpha} && \text{(inequality (22))} \\
 &\leq \min_{\mathcal{Z}: \kappa(\mathcal{Z}) \leq 1} f(\mathcal{W}_t) + \eta_t \langle s\mathcal{Z} - \mathcal{W}_t, \nabla\ell(\mathcal{W}_t) \rangle + \frac{L \left\| \frac{s}{\alpha}\mathcal{Z}_t - \mathcal{W}_t \right\|^2}{2} \eta_t^2 - \lambda\eta_t\kappa(\mathcal{W}_t) + \lambda\eta_t\frac{s}{\alpha} && \text{(definition of } \mathcal{Z}_t\text{)} \\
 &= \min_{\mathcal{Z}: \kappa(\mathcal{Z}) \leq s} f(\mathcal{W}_t) + \eta_t \langle \mathcal{Z} - \mathcal{W}_t, \nabla\ell(\mathcal{W}_t) \rangle + \frac{L \left\| \frac{s}{\alpha}\mathcal{Z}_t - \mathcal{W}_t \right\|^2}{2} \eta_t^2 - \lambda\eta_t\kappa(\mathcal{W}_t) + \lambda\eta_t\frac{s}{\alpha} && \text{(homogeneity of } \kappa\text{)} \\
 &\leq \min_{\mathcal{Z}: \kappa(\mathcal{Z}) \leq s} f(\mathcal{W}_t) + \eta_t(\ell(\mathcal{Z}) - \ell(\mathcal{W}_t)) + \frac{L \left\| \frac{s}{\alpha}\mathcal{Z}_t - \mathcal{W}_t \right\|^2}{2} \eta_t^2 - \lambda\eta_t \cdot \kappa(\mathcal{W}_t) + \lambda\eta_t\frac{s}{\alpha} && \text{(convexity of } \ell\text{)} \\
 &= (1 - \eta_t)f(\mathcal{W}_t) + \eta_t \min_{\mathcal{Z}: \kappa(\mathcal{Z}) \leq s} (\ell(\mathcal{Z}) + \lambda \cdot \frac{s}{\alpha}) + \frac{L \left\| \frac{s}{\alpha}\mathcal{Z}_t - \mathcal{W}_t \right\|^2}{2} \eta_t^2 \\
 &\leq (1 - \eta_t)f(\mathcal{W}_t) + \eta_t \frac{f(\mathcal{W})}{\alpha} + \frac{L \left\| \frac{s}{\alpha}\mathcal{Z}_t - \mathcal{W}_t \right\|^2}{2} \eta_t^2 && (\ell \geq 0 \text{ and } \alpha \in (0, 1]).
 \end{aligned}$$

Therefore,

$$f(\mathcal{W}_{t+1}) - \frac{f(\mathcal{W})}{\alpha} \leq (1 - \eta_t) \left( f(\mathcal{W}_t) - \frac{f(\mathcal{W})}{\alpha} \right) + \frac{L \left\| \frac{s}{\alpha}\mathcal{Z}_t - \mathcal{W}_t \right\|^2}{2} \eta_t^2.$$

Recall that  $\eta_t = \frac{2}{t+2}$ . An easy induction argument establishes that

$$f(\mathcal{W}_{t+1}) - \frac{f(\mathcal{W})}{\alpha} \leq \frac{2C}{t+3},$$

where  $C := \sup_t L \left\| \frac{s}{\alpha}\mathcal{Z}_t - \mathcal{W}_t \right\|^2 \leq 2L(\kappa^2(\mathcal{W})/\alpha^2 + D^2) < \infty$ , since  $\mathcal{W}_t$  is in the sublevel set of  $\{\mathcal{Z} : f(\mathcal{Z}) \leq f(\mathcal{W}_1)\}$ , whose radius is assumed to be bounded by  $D$ .  $\square$

**Theorem 2.** Fix  $\mathcal{A} \in \mathbb{R}^{d_1 \times \dots \times d_K}$  and  $t \geq \prod_{k=1}^K d_k$ , then

$$\|\mathcal{A}\|_{\text{tr}} = \min \left\{ \sum_{i=1}^t \|\mathbf{u}_i\|_2 \cdots \|\mathbf{z}_i\|_2 \right\} \quad (23)$$

$$= \min \left\{ \frac{1}{K} \sum_{i=1}^t \|\mathbf{u}_i\|_2^K + \cdots + \|\mathbf{z}_i\|_2^K \right\}, \quad (24)$$

where the minimum is taken w.r.t. all factorizations  $\mathcal{A} = \sum_{i=1}^t \mathbf{u}_i \otimes \cdots \otimes \mathbf{z}_i$ ,  $\mathbf{u}_i \in \mathbb{R}^{d_1}, \dots, \mathbf{z}_i \in \mathbb{R}^{d_K}$ .

*Proof.* We first note that the atomic set  $\mathbf{A}$  in (4) is compact, so is its convex hull  $\text{conv}(\mathbf{A})$ . Moreover,  $\mathbf{A}$  is connected. Recall that the trace norm is defined in (5) via the gauge function  $\kappa$ . Since  $\text{conv}(\mathbf{A})$  is compact with  $\mathbf{0}$  in its interior, we know the infimum in (5) is attained. Thus there exist  $\rho \geq 0$  and  $\mathcal{C} \in \text{conv}(\mathbf{A})$  so that  $\mathcal{A} = \rho\mathcal{C}$  and  $\|\mathcal{A}\|_{\text{tr}} = \kappa(\mathcal{A}) = \rho$ . Applying Caratheodory's theorem we know  $\mathcal{C} = \sum_{i=1}^t \sigma_i \hat{\mathbf{u}}_i \otimes \cdots \otimes \hat{\mathbf{z}}_i$  for some  $\sigma_i \geq 0$ ,  $\sum_i \sigma_i = 1$ ,  $\hat{\mathbf{u}}_i \otimes \cdots \otimes \hat{\mathbf{z}}_i \in \mathbf{A}$  and  $t \leq \prod_{k=1}^K d_k$ . Let  $\mathbf{u}_i = \sqrt[\kappa]{\rho\sigma_i} \hat{\mathbf{u}}_i, \dots, \mathbf{z}_i = \sqrt[\kappa]{\rho\sigma_i} \hat{\mathbf{z}}_i$  we know  $\|\mathcal{A}\|_{\text{tr}}$  is at least the right-hand side of (24).

On the other hand, for any  $\mathcal{A} = \sum_{i=1}^t \mathbf{u}_i \otimes \cdots \otimes \mathbf{z}_i$ , denoting  $\sigma_i = \|\mathbf{u}_i\|_2 \cdots \|\mathbf{z}_i\|_2$  ( $\neq 0$  w.l.o.g.), we have  $\mathcal{A} = (\sum_i \sigma_i) \cdot \sum_{i=1}^t \frac{\sigma_i}{\sum_{j=1}^t \sigma_j} \mathcal{A}_i$ , where  $\mathcal{A}_i := \frac{\mathbf{u}_i}{\|\mathbf{u}_i\|_2} \otimes \cdots \otimes \frac{\mathbf{z}_i}{\|\mathbf{z}_i\|_2} \in \mathbf{A}$ . Thus, appealing to the definition (5) we know  $\|\mathcal{A}\|_{\text{tr}} \leq \sum_i \sigma_i$ , i.e. (23) holds with  $\leq$ .

To complete the proof, we apply the arithmetic-geometric mean inequality:

$$\sum_i \|\mathbf{u}_i\|_2 \cdots \|\mathbf{z}_i\|_2 \leq \frac{1}{K} \sum_i \|\mathbf{u}_i\|_2^K + \cdots + \|\mathbf{z}_i\|_2^K.$$

(Of course, other elementary symmetric functions can be similarly used in Theorem 2.)  $\square$

## D Comparison with alternative variational forms

We compare the variational forms in Theorem 2 to some existing ones in the tensor literature. Firstly, the regularization function

$$\sum_i \|\mathbf{u}_i\|_2^2 + \cdots + \|\mathbf{z}_i\|_2^2, \quad (25)$$

is extensively used in finding CP decompositions, since otherwise the factors could blow up in a way that still maintains their sum, the so-called degeneracy problem [9]. A second reason for employing (25) is that it adds strict convexity w.r.t. each factor hence guarantees the convergence of block coordinate ascent. However, both reasons to promote (25), albeit valid, are weak; there are certainly other, perhaps even better, candidate regularizations. For instance, (24) enjoys both properties, with the additional equivalence to the trace norm, which potentially could lead to a low rank solution. The second variational form, appeared in [13], is

$$\|\mathcal{C}\|_{\mathbb{F}}^2 + \|U\|_{\mathbb{F}}^2 + \cdots + \|Z\|_{\mathbb{F}}^2, \quad (26)$$

where  $\mathcal{A} = \mathcal{C} \times_1 U \cdots \times_K Z$  is the Tucker decomposition. [13] used (26) to avoid the scaling ambiguity—a weak motivation for the particular form (26) indeed. Let us show that *neither* (25) *nor* (26) is equivalent to the trace norm. From this regard, the variational forms in Theorem 2 are advantageous and perhaps should be favored more often in practice.

**Example 3.** *We first prove (25) is not equivalent to the trace norm. Let  $\mathcal{A} = \sigma \mathbf{u} \otimes \cdots \otimes \mathbf{z}$  be a rank-1 tensor with  $\sigma > 0$  and  $\|\mathbf{u}\|_2 = \cdots = \|\mathbf{z}\|_2 = 1$ . It is easy to see that  $\|\mathcal{A}\|_{\text{tr}} = \sigma$ . Consider the function:*

$$f(\mathcal{A}) := \min \left\{ \sum_i \|\mathbf{u}_i\|_2^2 + \cdots + \|\mathbf{z}_i\|_2^2 \right\},$$

where the minimum is taken w.r.t. to all factorizations  $\mathcal{A} = \sum_i \mathbf{u}_i \otimes \cdots \otimes \mathbf{z}_i$ . Clearly,

$$f(\mathcal{A}) \leq K \cdot \sigma^{2/K}.$$

Choose an appropriate  $\sigma > 1$  we thus have  $f(\mathcal{A}) < \|\mathcal{A}\|_{\text{tr}}$ . Of course, for any positive constant  $c$ , we can choose appropriate  $\sigma$  such that  $c \cdot f(\mathcal{A}) < \|\mathcal{A}\|_{\text{tr}}$ . Thus (25) is not proportional to the trace norm.

For the function (26), we similarly define

$$g(\mathcal{A}) := \inf \left\{ \|\mathcal{C}\|_{\mathbb{F}}^2 + \|U\|_{\mathbb{F}}^2 + \cdots + \|Z\|_{\mathbb{F}}^2 \right\},$$

where the infimum is taken w.r.t. all Tucker decompositions  $\mathcal{A} = \mathcal{C} \times_1 U \cdots \times_K Z$ . This removes the dependence of (26) on a particular Tucker decomposition (which may not be unique). Consider the same rank-1 tensor  $\mathcal{A}$  as above, we have  $g(\mathcal{A}) \leq (K+1)\sigma^{2/(K+1)}$  while  $\|\mathcal{A}\|_{\text{tr}} = \sigma$ . Again, choosing  $\sigma$  large we have  $c \cdot g(\mathcal{A}) < \|\mathcal{A}\|_{\text{tr}}$  for any positive constant  $c$ . Note that even for  $K = 2$ ,  $g(\mathcal{A})$  is not proportional to the trace norm.

## E Efficient gradient computation

An efficient implementation for the surrogate problem (14) using state-of-the-art solvers (e.g. L-BFGS) relies on the efficient computation of the gradient of  $\ell$ . For simplicity, consider  $t = 1$  hence  $\mathcal{W} = \mathbf{u}_1 \otimes \cdots \otimes \mathbf{u}_K$ . The idea generalizes straightforwardly to all  $t$ . Let  $\mathcal{G} = \nabla \ell(\mathcal{W})$ . Using the chain rule,

$$\frac{\partial \ell}{\partial \mathbf{u}_k}(\mathcal{W}) = \mathcal{G} \times_1 \mathbf{u}_1^\top \cdots \times_{k-1} \mathbf{u}_{k-1}^\top \times_{k+1} \mathbf{u}_{k+1}^\top \cdots \times_K \mathbf{u}_K^\top.$$

The chain of product on the right-hand side costs  $O(\prod_k d_k)$ , thus a naive implementation for all factors would cost  $O(K \prod_k d_k)$ . Fortunately, using the identity in (17) we can reduce the cost by a factor of  $K$ . Define the forward and backward accumulators

$$\mathcal{F}_k := \mathcal{G} \times_1 \mathbf{u}_1^\top \dots \times_{k-1} \mathbf{u}_{k-1}^\top, \quad \mathcal{B}_k := \mathbf{u}_{k+1} \otimes \dots \otimes \mathbf{u}_K,$$

with  $\mathcal{F}_1 := \mathcal{G}$  and  $\mathcal{B}_K := 1$ . Then we have

$$\frac{\partial \ell}{\partial \mathbf{u}_k}(\mathcal{W}) = \mathcal{F}_k \times_{k+1} \mathbf{u}_{k+1}^\top \dots \times_K \mathbf{u}_K^\top = (\mathcal{F}_k)_{(k)} \mathcal{B}_k.$$

So we need only compute  $\{\mathcal{F}_k, \mathcal{B}_k\}$ , costing  $O(\prod_k d_k)$ . Since for all  $k$  the multiplication  $(\mathcal{F}_k)_{(k)} \mathcal{B}_k$  costs  $O(\prod_{j \geq k} d_j)$ , the overall time and space costs are both  $O(\prod_k d_k)$ . Clearly the computational savings are possible due to our explicit low-rank representation, which is not available in other matricization approaches.

## F Comparison on completing tensors with low-rank Tucker decompositions

We repeated all comparisons conducted for low CP rank on low Tucker rank:  $\mathcal{Z}^0 = \mathcal{S} \times_1 U_1 \times_2 U_2 \times_3 U_3$ , where  $\mathcal{S} \in \mathbb{R}^{r \times r \times r}$  and  $U_i \in \mathbb{R}^{n \times r}$ . Our setting here is exactly the same as Section 6.1. Again, all entries of  $\mathcal{S}$  and  $U_i$  were drawn i.i.d. from a unit normal. We set the default  $p=20$ ,  $\sigma=0.1$ ,  $r=5$  (hence CP rank  $\leq 25$ ), and  $n=50$ . Figure 7 shows that even in this case TTN still outperforms HaLRTC and RpLRTC (abbreviated as RP).

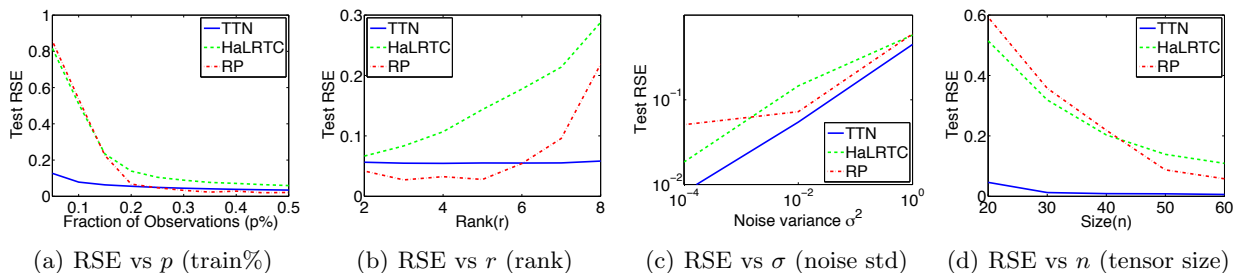


Figure 7: Test RSE on synthetic data generated by Tucker decomposition with rank  $r \times r \times r$ .

## G Robust PCA (RPCA)

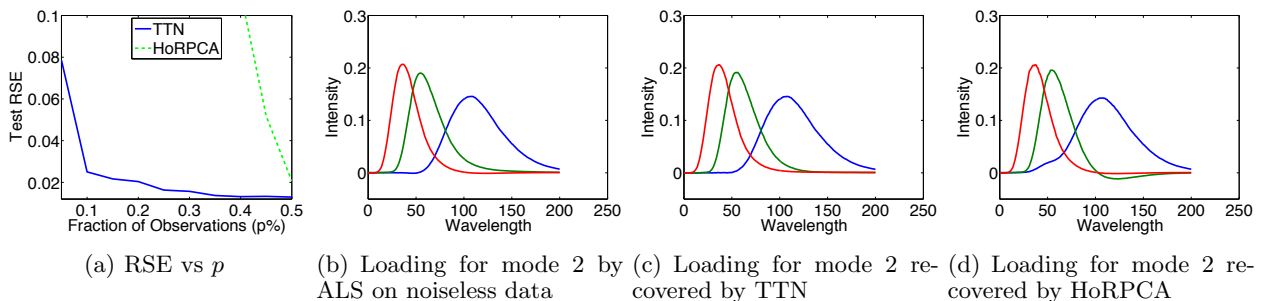


Figure 8: RPCA on amino acid data.  $p = 40$  in c-d.

Higher-order RPCA (HoRPCA), introduced in [16], solves  $\|\mathcal{W} - \mathcal{Z}\|_1 + \lambda \Omega(\mathcal{W})$ , where  $\Omega(\mathcal{W})$  is, say the sum of matrix trace norms of all mode unfoldings. The  $\ell_1$  loss used here aims to instill robustness to (large) outliers. To apply GCG to our proposed TTN as  $\Omega$ , we smoothed the  $\ell_1$  loss using a quadratic prox-function as in [39].

**Amino acid fluorescence data.** This dataset of 3-mode tensor records the fluorescence intensity of 5 solutions (mode 1), with 201 intervals of emission wavelength (mode 2) and 61 intervals of excitation wavelength (mode 3) [40]. Since the solutions contain three amino acids, the tensor has approximately low rank of 3. To mimic scatter outliers in chmometrics like [16], we randomly selected 10% entries on which i.i.d. corruptions uniformly distributed in  $[-100, 100]$  were added. Like above, we observe only a fraction,  $p\%$ , of the entries.

Figure 8(a) shows the RSE on the remaining  $(90 - p)\%$  entries, measured on the tensors recovered by using TTN and HoRPCA. Over a wide range of  $p$ , TTN achieves significantly lower RSE. For a more intuitive illustration,

we plot in Figure 8(b) the original loading (CP factor) for mode 2, as well as the loadings recovered by TTN and HoRPCA in Figures 8(c) and 8(d), respectively. Clearly TTN discovers a more faithful reconstruction of the CP factors in the presence of noise. Loadings for mode 3 are provided in Figure 9.

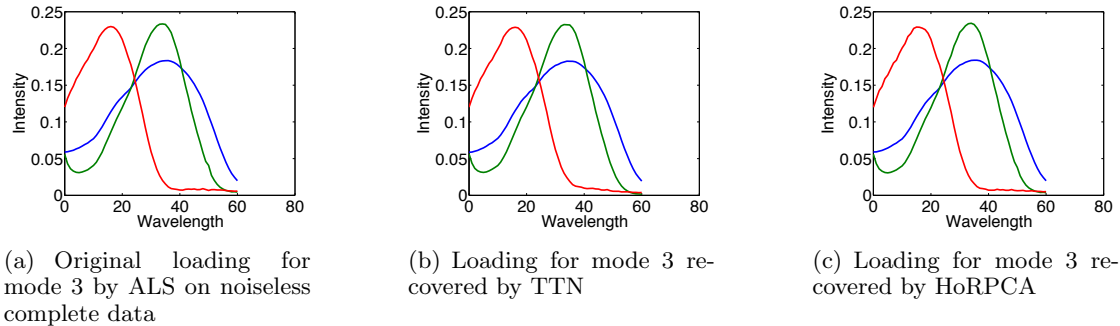


Figure 9: Loadings for mode 3 (excitation) with RPCA on amino acid data.  $p = 40$  in b-c.

We also experimented on another fluorescence dataset: Dorrit. It is a  $23 \times 116 \times 18$  tensor whose modes have the same meanings as in the amino acid data. However, it already contains scattering noise hence we no longer corrupted it. In Figure 10(c) and 10(d), we plot the loading for mode 2 recovered by TTN and HoRPCA, with comparison to that obtained from the noisy raw tensor (Figure 10(b)). All methods used the full tensor with no subsampling. Referring to the pure component loadings copied to Figure 10(a) from [41, Figure 2], one can observe that TTN attains a higher level of faithfulness than HoRPCA.

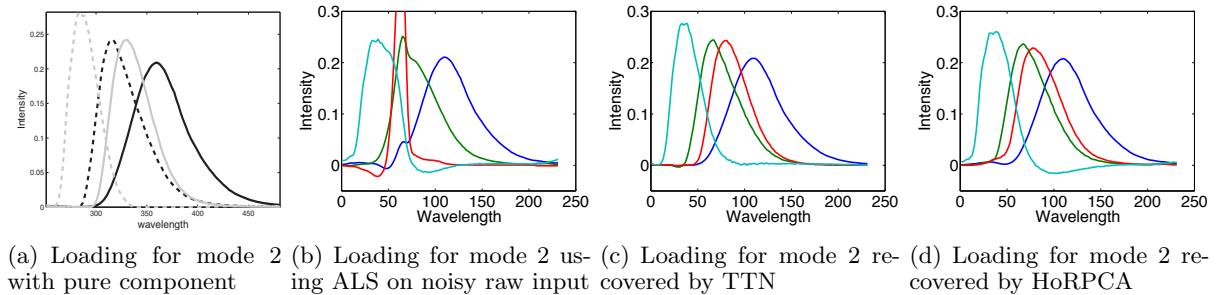


Figure 10: Loadings for mode 2 (emission) with RPCA on Dorrit data.  $p = 100$  in c-d.

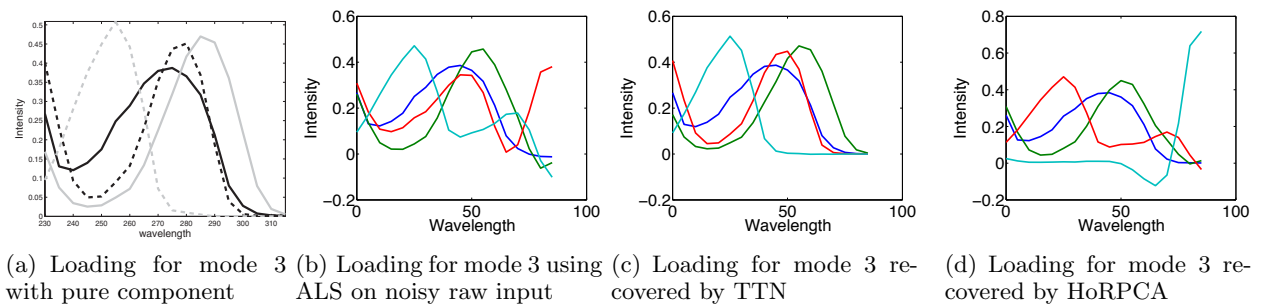


Figure 11: Loadings for mode 3 (excitation) with RPCA on Dorrit data.  $p = 100$  in c-d.