
On the extension of trace norm to tensors

Ryota Tomioka¹,

Kohei Hayashi²,

Hisashi Kashima¹

¹Department of Mathematical Informatics,
The University of Tokyo

²Graduate School of Information Science,
Nara Institute of Science and Technology

{tomioka, kashima}@mist.i.u-tokyo.ac.jp

kohei-h@is.naist.jp

Abstract

In this paper, we propose three extensions of trace norm for the minimization of tensor rank via convex optimization. One of the proposed extensions recovers partially observed tensor almost perfectly from a small fraction of observations.

1 Introduction

Higher order tensor decompositions have recently been studied intensively motivated by their usefulness in various fields including chemometrics, neuroimaging, and graph analysis [1]. Tensor decomposition methods aim to separate the factors spanning each modality of a given tensor and the interactions among the factors. The smaller the number of factors or the smaller the number of interactions, the more compact and succinct the decomposition is.

Matrix completion, or matrix estimation, has recently witnessed a great advance driven by powerful theory [2, 3] from compressed sensing and tools from convex optimization [4].

In this paper, we consider the problem of tensor completion, or more generally tensor estimation, which aims to recover an unknown tensor from partial (noisy) observations under the assumption that the underlying tensor admits a compact decomposition. This setting obviously includes the case of decomposing a fully observed tensor.

The aim of this paper is to extend the trace norm, which is the key component in matrix completion via convex optimization, for the tensor completion problem.

2 Low rank matrix and tensor

2.1 Rank of a matrix and the trace norm

The rank r of an $R \times C$ matrix \mathbf{X} is defined via the singular-value decomposition (SVD)

$$\mathbf{X} = \mathbf{U} \text{diag}(\sigma_1(\mathbf{X}), \sigma_2(\mathbf{X}), \dots, \sigma_r(\mathbf{X})) \mathbf{V}^\top,$$

where $\mathbf{U} \in \mathbb{R}^{R \times r}$ and $\mathbf{V} \in \mathbb{R}^{C \times r}$ are orthogonal matrices, and $\sigma_j(\mathbf{X})$ is the j th largest singular-value of \mathbf{X} . The matrix \mathbf{X} is called *low-rank* if the rank r is less than $\min(R, C)$. Unfortunately, the rank of a matrix is a nonconvex function, and the direct minimization of matrix rank is an NP-hard problem.

The trace norm is known to be the tightest convex lower bound of matrix rank [3] and is defined as the linear sum of singular values as follows:

$$\|\mathbf{X}\|_* = \sum_{j=1}^r \sigma_j(\mathbf{X}).$$

The trace norm allows us to estimate low-rank matrices via convex optimization with a theoretical guarantee [2]. Intuitively, the trace norm plays the role of the ℓ_1 -norm in the subset selection problem, for the estimation of low-rank matrices.

2.2 Rank of a tensor

We consider the k -rank of tensors, which is a direct generalization of the above definition of the matrix rank; see [1] for other definitions of tensor rank.

The k -rank of an K th-order tensor \mathcal{X} , denoted $\text{rank}_k(\mathcal{X})$, is defined as the rank of the mode- k unfolding $\mathbf{X}_{(k)}$ of \mathcal{X} . The tensor \mathcal{X} is called low-rank if any of its unfoldings is a low-rank matrix.

A rank- $(r_1, \dots, r_k, \dots, r_K)$ tensor \mathcal{X} of dimensions $n_1 \times \dots \times n_K$ can be written as

$$\mathcal{X} = \mathcal{G} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \cdots \times_K \mathbf{U}_K,$$

where \times_k denotes the k -mode matrix product, $\mathcal{G} \in \mathbb{R}^{r_1 \times \dots \times r_K}$ is called the *core tensor*, and $\mathbf{U}_k \in \mathbb{R}^{n_k \times r_k}$ ($n = 1, \dots, K$) are left singular-vectors from the SVD of the mode- k unfolding of \mathcal{X} . The above decomposition is called the Tucker decomposition [1].

Since the core tensor \mathcal{G} that corresponds to singular-values in the matrix case is not diagonal in general, it is not straightforward to generalize the trace norm from matrices to tensors.

3 Three strategies to extend the trace-norm regularization to tensors

In this section, we first consider a given tensor as a matrix and propose to minimize the trace norm of one of its unfoldings. Next, we extend this to the minimization of the weighted sum of the trace norms of the unfoldings. Finally, relaxing the condition that the tensor is *jointly* low-rank in every mode in the second approach, we propose a mixture approach.

3.1 Tensor as a matrix

The definition of a low-rank tensor in the previous section implies that a low-rank tensor *is* a low-rank matrix when unfolded appropriately.

Therefore, for the reconstruction of partly observed tensor, we can solve the following problem:

$$\underset{\mathcal{X} \in \mathbb{R}^{n_1 \times \dots \times n_K}}{\text{minimize}} \quad \frac{1}{2\lambda} \|\Omega(\mathcal{X}) - \mathbf{y}\|^2 + \|\mathbf{X}_{(k)}\|_*, \quad (1)$$

where $\mathbf{X}_{(k)}$ is the mode- k unfolding of \mathcal{X} , $\mathbf{y} \in \mathbb{R}^M$ is the vector of observations, and $\Omega : \mathbb{R}^{n_1 \times \dots \times n_K} \rightarrow \mathbb{R}^M$ is a linear operator that reshapes the prespecified (possibly overlapping) elements of the input tensor into an M dimensional vector; M is the number of observations.

Since the estimation procedure (1) is essentially an estimation of a low-rank matrix $\mathbf{X}_{(k)}$, we know that $O(\tilde{n}_k^{6/5} r_k \log(\tilde{n}_k))$ samples are enough to perfectly recover the unknown true tensor \mathcal{X}^* , where $r_k = \text{rank}_k(\mathcal{X}^*)$ and $\tilde{n}_k = \max(n_k, \prod_{k' \neq k} n_{k'})$, if the rank r_k is not too high [2].

Note that when we can estimate the mode- k unfolding of \mathcal{X}^* perfectly, we can also recover the whole \mathcal{X}^* perfectly, including the ranks of the modes we did not use during the estimation.

However, the success of the above procedure is conditioned on the choice of the mode to unfold the tensor. If we choose a mode with a large rank, even if there are other modes with smaller ranks, we cannot hope to recover the tensor from a small number of samples.

3.2 Constrained optimization of low rank tensors

In order to exploit the rank deficiency of more than one mode, it is natural to consider the following extension of the estimation procedure (1)

$$\underset{\mathcal{X} \in \mathbb{R}^{n_1 \times \dots \times n_K}}{\text{minimize}} \quad \frac{1}{2\lambda} \|\Omega(\mathcal{X}) - \mathbf{y}\|^2 + \sum_{k=1}^K \gamma_k \|\mathbf{X}_{(k)}\|_*$$

This is a convex optimization problem, because it can be reformulated as follows:

$$\underset{\mathbf{x}, \mathbf{Z}_1, \dots, \mathbf{Z}_K}{\text{minimize}} \quad \frac{1}{2\lambda} \|\Omega \mathbf{x} - \mathbf{y}\|^2 + \sum_{k=1}^K \gamma_k \|\mathbf{Z}_k\|_*, \quad (2)$$

$$\text{subject to} \quad \mathbf{P}_k \mathbf{x} = \mathbf{z}_k \quad (k = 1, \dots, K), \quad (3)$$

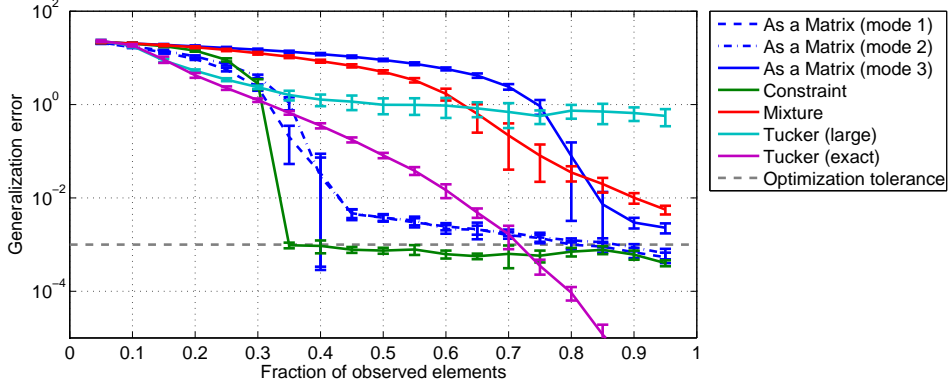


Figure 1: Comparison of three strategies, tensor as a matrix (“As a Matrix”), constrained optimization (“Constraint”), and mixture of low-rank tensors (“Mixture”). Also the Tucker decomposition with 20% higher rank (“large”) and with the correct rank (“exact”) implemented in the N-way toolbox [6] are included as baselines. The generalization error is plotted against the fraction of observed elements of the underlying low-rank tensor. Also the tolerance of optimization (10^{-3}) is shown.

where $\mathbf{x} \in \mathbb{R}^N$ is the vectorization of \mathcal{X} ($N = \prod_{k=1}^K n_k$), \mathbf{P}_k is the matrix representation of mode- k unfolding (note that \mathbf{P}_k is a permutation matrix; thus $\mathbf{P}_k^\top \mathbf{P}_k = \mathbf{I}_N$), $\mathbf{Z}_k \in \mathbb{R}^{n_k \times N/n_k}$ is a matrix of the same size as the mode- k unfolding of \mathcal{X} , and \mathbf{z}_k is the vectorization of \mathbf{Z}_k . With a slight abuse of notation $\mathbf{\Omega} \in \mathbb{R}^{M \times N}$ denotes the observation operator as a matrix.

This approach was considered earlier in [5], but they relaxed the constraints (3) into penalty terms, which is more similar to the approach we discuss in the next subsection.

3.3 Mixture of low-rank tensors

The optimization problem (2) is a constrained optimization problem and requires sophisticated tools to optimize. One simpler alternative is to predict with a mixture of K tensors; each mixture component is regularized by the trace norm to be low-rank in each mode. More specifically, we solve the following minimization problem:

$$\underset{\mathbf{Z}_1, \dots, \mathbf{Z}_K}{\text{minimize}} \quad \frac{1}{2\lambda} \left\| \mathbf{\Omega} \left(\frac{1}{K} \sum_{k=1}^K \mathbf{P}_k^\top \mathbf{z}_k \right) - \mathbf{y} \right\|^2 + \sum_{k=1}^K \gamma_k \|\mathbf{Z}_k\|_* \quad (4)$$

Note that when $\mathbf{x} = \mathbf{P}_k^\top \mathbf{z}_k$ for all $k = 1, \dots, K$, the problem (4) reduces to the problem (2).

4 Numerical experiments

We randomly generated a rank-(7,8,9) tensor of dimensions (50,50,20) by drawing the core from the standard normal distribution and multiplying its each mode by an orthonormal factor randomly drawn from the Haar measure. We randomly selected some elements of the true tensor for training and kept the remaining elements for testing. Alternating direction method of multiplier [7] with tolerance 10^{-3} was used for the optimization problems (1), (2), and (4). We choose $\gamma_k = 1$ for simplicity in the later two approaches. For the first two approaches, $\lambda \rightarrow 0$ (zero training error) was used. For the third approach (4), $\lambda = 0.001$ was used. Computational details are omitted for the sake of brevity. The Tucker decomposition algorithm `tucker` in the N-way toolbox [6] is also included as a baseline, for which we used the correct rank (“exact”) and the 20% higher rank (“large”). Note that all proposed approaches can find the rank automatically. The generalization error is defined as the square-root of the sum of squared differences between the true and the estimated tensors over the unobserved elements. For the “As a Matrix” strategy, error for each mode is reported. The experiment was repeated 10 times and averaged.

Figure 1 shows the result of tensor completion using three strategies we proposed above, as well as the Tucker decomposition. The proposed “Constraint” approach shows a sharp threshold behaviour

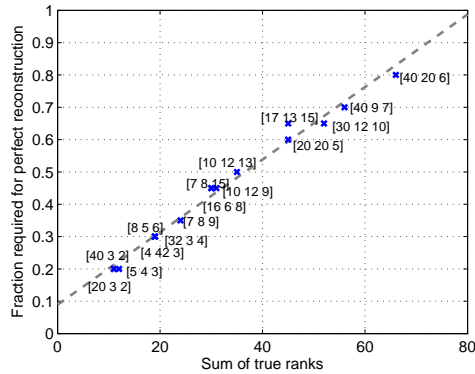


Figure 2: Fraction of observations at the threshold plotted against the sum of true ranks. Numbers in the brackets denote the k -rank of the underlying tensor. The dimension of the tensor is (50,50,20).

around 35% observation from very bad fit (generalization error > 1) to almost perfect fit (generalization error $\simeq 10^{-3}$). The “As a Matrix” approach also show similar transition for mode 1 and mode 2 (around 40%), and mode 3 (around 80%), but even the first transition is slower than the “Constraint” approach. The “Mixture” approach shows no clear transition. Tucker shows the fastest decrease in the generalization error, but when the rank is misspecified (“large”), the error remains almost constant; even when the correct rank is known (“exact”), the convergence is slower than the proposed “Constraint” approach.

We have further investigated the condition for the threshold behaviour using the proposed “Constraint” approach. In Figure 2, we can see that the fraction of observations required to perfectly recover an unknown tensor is roughly proportional to the sum of the rank of the underlying tensor, where we define the reconstruction to be perfect when the mean generalization error is less than 0.01.

5 Summary

In this paper we have proposed three strategies to extend the trace norm to tensor rank minimization, and we have compared them on a simulated tensor completion problem. We have found that tensor completion using the “Constraint” approach shows nearly perfect reconstruction from only 35% observations. There is no need to specify the rank of the decomposition as in the conventional Tucker decomposition approach. The proposed approach shows a sharp threshold behaviour and we have found that the fraction of samples at the threshold is roughly proportional to the sum of ranks of the underlying tensor. Further analysis is necessary to explain the threshold behaviour.

References

- [1] T. G. Kolda and B. W. Bader. Tensor decompositions and applications. *SIAM Review*, 51(3):455–500, 2009.
- [2] E. J. Candes and B. Recht. Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9(6):717–772, 2009.
- [3] B. Recht, M. Fazel, and P.A. Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Review*, 52(3):471–501, 2010.
- [4] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [5] J. Liu, P. Musialski, P. Wonka, and J. Ye. Tensor completion for estimating missing values in visual data. In *Prof. ICCV*, 2009.
- [6] C. A. Andersson and R. Bro. The n-way toolbox for matlab. *Chemometrics & Intelligent Laboratory Systems*, 52(1):1–4, 2000. <http://www.models.life.ku.dk/source/nwaytoolbox/>.
- [7] D. Gabay and B. Mercier. A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Comput. Math. Appl.*, 2(1):17–40, 1976.