



Rank Properties and Computational Methods for Orthogonal Tensor Decompositions

Chao Zeng¹

Received: 4 June 2022 / Revised: 23 October 2022 / Accepted: 31 October 2022 /
Published online: 23 November 2022

© The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2022

Abstract

The orthogonal decomposition factorizes a tensor into a sum of an orthogonal list of rank-one tensors. The corresponding rank is called orthogonal rank. We present several properties of orthogonal rank, which are different from those of tensor rank in many aspects. For instance, a subtensor may have a larger orthogonal rank than the whole tensor. To fit the orthogonal decomposition, we propose an algorithm based on the augmented Lagrangian method. The gradient of the objective function has a nice structure, inspiring us to use gradient-based optimization methods to solve it. We guarantee the orthogonality by a novel orthogonalization process. Numerical experiments show that the proposed method has a great advantage over the existing methods for strongly orthogonal decompositions in terms of the approximation error.

Keywords Orthogonal tensor decomposition · Orthogonal rank · Augmented Lagrangian method · L-BFGS · Orthogonalization

Mathematics Subject Classification 15A69 · 49M27 · 90C26 · 90C30

1 Introduction

Given a tensor $\mathcal{A} \in \mathbb{R}^{I_1 \times \dots \times I_N}$, the CANDECOMP/PARAFAC (CP) decomposition factorizes it into a sum of rank-one tensors:

$$\mathcal{A} = \sum_{k=1}^K \mathbf{v}_k^{(1)} \otimes \dots \otimes \mathbf{v}_k^{(N)},$$

where $\mathbf{v}_k^{(n)} \in \mathbb{R}^{I_n}$, $k = 1, \dots, K$, $n = 1, \dots, N$. Usually, it is difficult to determine the number K for expressing \mathcal{A} exactly [15, 16]. Hence, the following approximate CP decomposition

✉ Chao Zeng
zengchao@nankai.edu.cn

¹ School of Mathematical Sciences and LPMC, Nankai University, Tianjin, China

is more meaningful in practical applications:

$$\min_{\mathbf{v}_r^{(n)} \in \mathbb{R}^{I_n}} \left\| \mathcal{A} - \sum_{r=1}^R \mathbf{v}_r^{(1)} \otimes \cdots \otimes \mathbf{v}_r^{(N)} \right\|,$$

where $\|\cdot\|$ is the Frobenius norm and R is a prescribed number. This problem is just to find a best rank- R approximation to \mathcal{A} . Unfortunately, this problem has no solution in general [9, 20].

As mentioned in Ref. [9], the major open question in tensor approximation is how to overcome the ill-posedness of the low rank approximation problem. One natural strategy is to impose orthogonality constraints, because the orthogonality is an inherent property of second-order tensor rank decompositions, i.e., matrix singular value decompositions (SVD). The orthogonal tensor decomposition can be traced back to [6] for the symmetric case, and then is studied in [17] for the general case:

$$\mathcal{A} = \sum_{r=1}^R \mathcal{T}_r, \quad \text{where } \text{rank}(\mathcal{T}_r) = 1 \text{ and } \langle \mathcal{T}_s, \mathcal{T}_t \rangle = 0 \text{ for all } 1 \leq s \neq t \leq R, \quad (1)$$

where $\langle \cdot, \cdot \rangle$ is the inner product that induces the Frobenius norm; see Sect. 2.1 for details. In [23], the orthogonality constraint is extended to general angular constraints, where several properties including the existence, uniqueness and exact recoverability are discussed. As a special case of decompositions with angular constraints, the orthogonal tensor decomposition also has these properties.

The greedy approach presented in Ref. [17] is the earliest method for computing a low orthogonal rank approximation to a given tensor, where one rank-one component is updated in one iteration. Specifically, suppose we have obtained k rank-one components. The $(k + 1)$ st rank-one component is updated by

$$\begin{aligned} \min_{\mathcal{U}} \quad & \left\| \mathcal{A} - \sum_{r=1}^k \mathcal{T}_r - \mathcal{U} \right\| \\ \text{s.t.} \quad & \text{rank}(\mathcal{U}) = 1 \quad \text{and} \quad \langle \mathcal{T}_r, \mathcal{U} \rangle = 0, \quad r = 1, \dots, k. \end{aligned}$$

This method is reasonable only if the Eckart-Young theorem [11] can be extended to the orthogonal decomposition, i.e., the best low orthogonal rank approximation can be obtained by truncating the orthogonal rank decomposition (see Sect. 3 for the definition). Refer to [17, Section 5] for details. However, a counterexample presented in Ref. [18] shows that such an extension is not possible. Suppose $\mathcal{T}_r = \otimes_{n=1}^N \mathbf{v}_r^{(n)}$ in (1). The orthogonality constraint has the following form

$$\prod_{n=1}^N \langle \mathbf{v}_s^{(n)}, \mathbf{v}_t^{(n)} \rangle = 0 \quad \text{for all } s \neq t.$$

This means that there exists at least one $m \in \{1, \dots, N\}$ such that $\langle \mathbf{v}_s^{(m)}, \mathbf{v}_t^{(m)} \rangle = 0$. However, we cannot determine the number m for different pairs of s, t . This is the main difficulty in fitting orthogonal decompositions. Practical existing algorithms are proposed by fixing the number m ; see [5, 13, 31, 34, 35]. Actually, these algorithms are aimed at strongly orthogonal decompositions, whose one or more factor matrices are orthogonal. All these algorithms follow a similar framework, by combining the alternating minimization method and the polar decomposition. For factor matrices with general angular constraints, a proximal gradient

algorithm is proposed in Ref. [26]. In [24], the Jacobi SVD algorithm is extended to reduce a tensor to a form with the ℓ_2 norm of the diagonal vector being maximized. The resulting form is not diagonal and hence this is not an algorithm for orthogonal decompositions discussed in this work.

In this paper, we first study orthogonal rank. We find that there are many differences between orthogonal rank and tensor rank. Orthogonal rank may be variant under the invertible n -mode product, a subtensor may have a larger orthogonal rank than the whole tensor, and orthogonal rank is lower semicontinuous. A refined upper bound of orthogonal rank [22] is presented. As for the algorithm, we employ the augmented Lagrangian method to convert (1) into an unconstrained problem. We find that the gradient of the objective function has a good structure. Therefore, we use gradient-based optimization methods to solve each subproblem. To guarantee the orthogonality of the final result, we develop an orthogonalization process. Numerical experiments show that our method has a great advantage over the existing methods for strongly orthogonal decompositions in terms of the approximation error.

The rest of this paper is organized as follows. Section 2 recalls some preliminary materials. In Sect. 3, we present several properties of orthogonal rank. The algorithm is proposed in Sect. 4. Experimental results are given in Sect. 5. Conclusions are presented in Sect. 6.

Notation

We use bold-face lowercase letters ($\mathbf{a}, \mathbf{b}, \dots$) to denote vectors, bold-face capitals ($\mathbf{A}, \mathbf{B}, \dots$) to denote matrices and calligraphic letters ($\mathcal{A}, \mathcal{B}, \dots$) to denote tensors. The notations \mathbf{I} and $\mathbf{0}$ denote the identity matrix and the zero matrix of suitable dimensions, respectively. The (i_1, i_2, \dots, i_N) th element of \mathcal{A} is denoted by $a_{i_1 i_2 \dots i_N}$. The n -mode product of a tensor \mathcal{A} by a matrix \mathbf{M} is denoted by $\mathbf{M} \cdot_n \mathcal{A}$. Following [9], we write $\mathbf{M}_1 \cdot_1 \dots \mathbf{M}_N \cdot_N \mathcal{A}$ more concisely as $(\mathbf{M}_1, \dots, \mathbf{M}_N) \cdot \mathcal{A}$. The mode- n unfolding matrix is denoted by $\mathbf{A}_{(n)}$, whose columns are all mode- n fibers of \mathcal{A} , $n = 1, \dots, N$.

2 Preliminaries

2.1 Inner Product, Angle and Orthogonality

Let $\mathcal{A}, \mathcal{B} \in \mathbb{R}^{I_1 \times \dots \times I_N}$. The inner product of \mathcal{A}, \mathcal{B} is defined by

$$\langle \mathcal{A}, \mathcal{B} \rangle := \sum_{i_1=1}^{I_1} \dots \sum_{i_N=1}^{I_N} a_{i_1, \dots, i_N} b_{i_1, \dots, i_N},$$

and the Frobenius norm of \mathcal{A} induced by this inner product is $\|\mathcal{A}\| = \sqrt{\langle \mathcal{A}, \mathcal{A} \rangle}$. Let $\mathcal{U} = \mathbf{u}^{(1)} \otimes \dots \otimes \mathbf{u}^{(N)}$ and $\mathcal{V} = \mathbf{v}^{(1)} \otimes \dots \otimes \mathbf{v}^{(N)}$. Then

$$\langle \mathcal{U}, \mathcal{V} \rangle = \prod_{n=1}^N \langle \mathbf{u}^{(n)}, \mathbf{v}^{(n)} \rangle \quad \text{and} \quad \|\mathcal{U}\| = \prod_{n=1}^N \|\mathbf{u}^{(n)}\|. \tag{2}$$

We say that \mathcal{A} is a *unit* tensor if $\|\mathcal{A}\| = 1$.

The *angle* between \mathcal{A}, \mathcal{B} is defined by

$$\angle(\mathcal{A}, \mathcal{B}) := \arccos \left\langle \frac{\mathcal{A}}{\|\mathcal{A}\|}, \frac{\mathcal{B}}{\|\mathcal{B}\|} \right\rangle. \tag{3}$$

Two tensors \mathcal{A}, \mathcal{B} are *orthogonal* ($\mathcal{A} \perp \mathcal{B}$) if $\langle \mathcal{A}, \mathcal{B} \rangle = 0$, i.e., $\angle(\mathcal{A}, \mathcal{B}) = \frac{\pi}{2}$. In (2), \mathcal{U} and \mathcal{V} are orthogonal if $\prod_{n=1}^N \langle \mathbf{u}^{(n)}, \mathbf{v}^{(n)} \rangle = 0$. This leads to other options for defining orthogonality of two rank-one tensors. Given $1 \leq i_1 < \dots < i_M \leq N$, we say that \mathcal{U} and \mathcal{V} are (i_1, \dots, i_M) -orthogonal if

$$\langle \mathbf{u}^{(i_m)}, \mathbf{v}^{(i_m)} \rangle = 0 \quad \forall 1 \leq m \leq M.$$

If $M = N$, we say that \mathcal{U} and \mathcal{V} are *completely orthogonal*.

A list of tensors $\mathcal{T}_1, \dots, \mathcal{T}_m$ is said to be orthogonal if $\langle \mathcal{T}_i, \mathcal{T}_j \rangle = 0$ for all distinct $i, j \in \{1, \dots, m\}$. An orthogonal list of tensors is an orthonormal list if each of its elements is a unit tensor. Similarly, we can define an (i_1, \dots, i_M) -orthogonal list of rank-one tensors.

2.2 CP Decompositions and Tensor Rank

The CP decomposition factorizes a tensor into a sum of rank-one tensors:

$$\mathcal{A} = \sum_{r=1}^R \mathbf{v}_r^{(1)} \otimes \dots \otimes \mathbf{v}_r^{(N)} := \llbracket \mathbf{V}^{(1)}, \dots, \mathbf{V}^{(N)} \rrbracket, \tag{4}$$

where the n th factor matrix is

$$\mathbf{V}^{(n)} = \begin{bmatrix} \mathbf{v}_1^{(n)} & \dots & \mathbf{v}_R^{(n)} \end{bmatrix}.$$

An interesting property of tensors is that their CP decompositions are often unique. Refer to [19, Section 3.2] for detailed introductions. The most famous results [21, 30] on the uniqueness condition depend on the concept of k -rank. The k -rank of a matrix \mathbf{M} , denoted by $k_{\mathbf{M}}$, is the largest integer such that every set containing $k_{\mathbf{M}}$ columns of \mathbf{M} is linearly independent. For the CP decomposition (4), its uniqueness condition presented in [30] is

$$\sum_{n=1}^N k_{\mathbf{V}^{(n)}} \geq 2R + N - 1. \tag{5}$$

Clearly, if (4) is unique, we must have $R = \text{rank}(\mathcal{A})$.

The rank of \mathcal{A} is defined by $\text{rank}(\mathcal{A}) := \min \left\{ R : \mathcal{A} = \sum_{r=1}^R \mathbf{v}_r^{(1)} \otimes \dots \otimes \mathbf{v}_r^{(N)} \right\}$. Given $R > 0$, the following problem

$$\min_{\text{rank}(\mathcal{B}) \leq R} \|\mathcal{A} - \mathcal{B}\| \tag{6}$$

aims to find the *best rank- R approximation* of \mathcal{A} . However, (6) has no solution in general [9, 20]. This is due to the following feature of tensor rank.

Proposition 2.1 ([9]) *Let $R \geq 2$. The set $\{\mathcal{A} \in \mathbb{R}^{I_1 \times \dots \times I_N} : \text{rank}(\mathcal{A}) \leq R\}$ is not closed in the normed space $\mathbb{R}^{I_1 \times \dots \times I_N}$. That is, the function $\text{rank}(\mathcal{A})$ is not lower semicontinuous.*

2.3 Orthogonal Decompositions

The orthogonal decomposition factorizes a tensor into a sum of an orthogonal list of rank-one tensors:

$$\mathcal{A} = \sum_{r=1}^R \mathcal{T}_r, \quad \text{where } \text{rank}(\mathcal{T}_r) = 1 \text{ and } \mathcal{T}_s \perp \mathcal{T}_t \text{ for all } 1 \leq s \neq t \leq R. \tag{7}$$

The following lemma can be obtained by a direct calculation based on (2).

Lemma 2.2 *The decomposition (4) is an orthogonal decomposition if and only if $\mathbf{V}^{(1)\top} \mathbf{V}^{(1)} \otimes \dots \otimes \mathbf{V}^{(N)\top} \mathbf{V}^{(N)}$ is diagonal, where “ \otimes ” is the Hadamard product.*

The (i_1, \dots, i_M) -orthogonal decomposition factorizes a tensor into a sum of an (i_1, \dots, i_M) -orthogonal list of rank-one tensors. Any type of (i_1, \dots, i_M) -orthogonal decomposition is called *strongly orthogonal decomposition*¹. Clearly, a strongly orthogonal decomposition is also an orthogonal decomposition. However, we are not in general guaranteed that a strongly orthogonal decomposition exists. Simple examples include the tensors with $\text{rank}(\mathcal{A}) > \max\{I_1, \dots, I_N\}$ ². This is because an (i_1, \dots, i_M) -orthogonal list consists of at most $\min\{I_{i_1}, \dots, I_{i_M}\}$ elements. Related discussions can be found in Ref. [5, 17].

There is a lot of research on strongly orthogonal decompositions. The $(1, \dots, N)$ -orthogonal decomposition, also called the *completely orthogonal decomposition*, is discussed in Ref. [5]. The (n) -orthogonality, where $1 \leq n \leq N$, is considered in Ref. [31, 34]. General strongly orthogonal decompositions are considered in Ref. [13, 35].

3 Properties of Orthogonal Rank

The *orthogonal rank* of \mathcal{A} is defined by

$$\text{rank}_\perp(\mathcal{A}) := \min \left\{ R \in \mathbb{N} : \mathcal{A} = \sum_{r=1}^R \mathcal{T}_r, \text{rank}(\mathcal{T}_r) = 1, \mathcal{T}_s \perp \mathcal{T}_t \text{ for all } 1 \leq s \neq t \leq R \right\}.$$

If $R = \text{rank}_\perp(\mathcal{A})$ in (7), then (7) is called an *orthogonal rank decomposition*.

Clearly, $\text{rank}_\perp(\mathcal{A}) \geq \text{rank}(\mathcal{A})$. The following lemma gives a necessary and sufficient condition for $\text{rank}_\perp(\mathcal{A}) = \text{rank}(\mathcal{A})$ under some assumptions.

Lemma 3.1 *Let $R \geq 2$, $\mathbf{V}^{(n)} \in \mathbb{R}^{I_n \times R}$ for $n = 1, \dots, N$ and $\mathcal{A} = \llbracket \mathbf{V}^{(1)}, \dots, \mathbf{V}^{(N)} \rrbracket$. Suppose $\text{rank}(\mathbf{V}^{(n)}) = R$ for all $n = 1, \dots, N$. Then $\text{rank}(\mathcal{A}) = \text{rank}_\perp(\mathcal{A}) = R$ if and only if $\mathbf{V}^{(1)\top} \mathbf{V}^{(1)} \otimes \dots \otimes \mathbf{V}^{(N)\top} \mathbf{V}^{(N)}$ is diagonal.*

Proof Since $\text{rank}(\mathbf{V}^{(n)}) = R$ and $R \geq 2$, we have

$$\sum_{n=1}^N k_{\mathbf{V}^{(n)}} = NR \geq 2R + N - 1.$$

By (5), this decomposition is unique and $\text{rank}(\mathcal{A}) = R$.

By Lemma 2.2, if $\mathbf{V}^{(1)\top} \mathbf{V}^{(1)} \otimes \dots \otimes \mathbf{V}^{(N)\top} \mathbf{V}^{(N)}$ is diagonal, then this decomposition is an orthogonal decomposition. Hence,

$$\text{rank}(\mathcal{A}) \leq \text{rank}_\perp(\mathcal{A}) \leq R = \text{rank}(\mathcal{A}) \Rightarrow \text{rank}(\mathcal{A}) = \text{rank}_\perp(\mathcal{A}) = R.$$

Also by Lemma 2.2, if $\mathbf{V}^{(1)\top} \mathbf{V}^{(1)} \otimes \dots \otimes \mathbf{V}^{(N)\top} \mathbf{V}^{(N)}$ is not diagonal, then this decomposition is not an orthogonal decomposition. Due to the uniqueness, there does not exist an orthogonal decomposition with R terms, i.e., $\text{rank}_\perp(\mathcal{A}) > R$. □

¹ Strongly orthogonal decomposition has a different definition in Ref. [17].

² Such tensors exist. See [9, Lemma 4.7] for an example.

Suppose \mathcal{A} is a subtensor of \mathcal{B} , then $\text{rank}(\mathcal{A}) \leq \text{rank}(\mathcal{B})$. It comes as a surprise that the analog does not hold for orthogonal rank. See the next proposition.

Proposition 3.2 *Let $R \geq 2, \mathbf{V}^{(n)} \in \mathbb{R}^{I_n \times R}$ for $n = 1, \dots, N$ and $\mathcal{A} = \llbracket \mathbf{V}^{(1)}, \dots, \mathbf{V}^{(N)} \rrbracket$. Suppose $\text{rank}(\mathbf{V}^{(n)}) = R$ for all $n = 1, \dots, N$. If $\mathbf{V}^{(1)\top} \mathbf{V}^{(1)} \otimes \dots \otimes \mathbf{V}^{(N)\top} \mathbf{V}^{(N)}$ is not diagonal, then there exists a tensor \mathcal{B} such that*

$$\mathcal{A} \text{ is a subtensor of } \mathcal{B} \text{ and } \text{rank}_{\perp}(\mathcal{B}) < \text{rank}_{\perp}(\mathcal{A}).$$

Proof We can find a sufficiently large t such that $t\mathbf{I} - \mathbf{V}^{(1)\top} \mathbf{V}^{(1)}$ is positive semidefinite. Then there exists a matrix \mathbf{M} with R columns such that

$$t\mathbf{I} - \mathbf{V}^{(1)\top} \mathbf{V}^{(1)} = \mathbf{M}^{\top} \mathbf{M}.$$

Define $\mathbf{V} = \begin{bmatrix} \mathbf{V}^{(1)} \\ \mathbf{M} \end{bmatrix}$. Then $\mathcal{B} = \llbracket \mathbf{V}, \mathbf{V}^{(2)}, \dots, \mathbf{V}^{(N)} \rrbracket$ is an orthogonal decomposition. By Lemma 3.1, we have $\text{rank}_{\perp}(\mathcal{B}) = R < \text{rank}_{\perp}(\mathcal{A})$. □

Remark 3.3 If we regard a matrix as a two-dimensional dataset, the value of rank represents the real quantity of information that the matrix embodies. That is, a low-rank matrix embodies a lot of redundant information. Therefore, it makes sense that the rank of a submatrix is smaller than that of the whole matrix. There are various ways to decompose a tensor into a sum of rank-one tensors, and each of these decompositions leads naturally to a concept of tensor rank. Proposition 3.2 shows the flaw of orthogonal rank in representing the real quantity of information that a tensor embodies.

A basic property of tensor rank is its invariance under the invertible n -mode product. If \mathbf{M}_n is invertible for $n = 1, \dots, N$, [9, Lemma 2.3] tells us that

$$\text{rank}((\mathbf{M}_1, \dots, \mathbf{M}_N) \cdot \mathcal{A}) = \text{rank}(\mathcal{A}).$$

However, the analog does not hold for orthogonal rank. Counterexamples can be constructed based on Lemma 3.1. Due to the fact that $\text{rank}(\mathbf{V}^{(1)}) = R$, there exists an invertible matrix $\mathbf{M} \in \mathbb{R}^{I_1 \times I_1}$ satisfying $\mathbf{M}(:, 1 : R) = \mathbf{V}^{(1)}$. Then $\mathbf{M}^{-1} \mathbf{V}^{(1)} = \begin{bmatrix} \mathbf{I} \\ \mathbf{0} \end{bmatrix}$ and $\mathbf{M}^{-1} \cdot_1 \mathcal{A} = \llbracket \mathbf{M}^{-1} \mathbf{V}^{(1)}, \mathbf{V}^{(2)}, \dots, \mathbf{V}^{(N)} \rrbracket$ is an orthogonal decomposition. Therefore,

$$\text{rank}_{\perp}(\mathbf{M}^{-1} \cdot_1 \mathcal{A}) = \text{rank}(\mathbf{M}^{-1} \cdot_1 \mathcal{A}) = \text{rank}(\mathcal{A}) < \text{rank}_{\perp}(\mathcal{A}).$$

If the n -mode product is orthogonal, we have the following lemma.

Lemma 3.4 *Let $\mathcal{A} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ and $\mathbf{M}_n \in \mathbb{R}^{I_n \times I_n}$ be orthogonal for $n = 1, \dots, N$. Then*

$$\text{rank}_{\perp}((\mathbf{M}_1, \dots, \mathbf{M}_N) \cdot \mathcal{A}) = \text{rank}_{\perp}(\mathcal{A}).$$

Proof Suppose $\mathcal{A} = \llbracket \mathbf{V}^{(1)}, \dots, \mathbf{V}^{(N)} \rrbracket$ is an orthogonal decomposition. Then $(\mathbf{M}_1, \dots, \mathbf{M}_N) \cdot \mathcal{A} = \llbracket \mathbf{M}_1 \mathbf{V}^{(1)}, \dots, \mathbf{M}_N \mathbf{V}^{(N)} \rrbracket$ and $(\mathbf{M}_1 \mathbf{V}^{(1)})^{\top} \mathbf{M}_1 \mathbf{V}^{(1)} \otimes \dots \otimes (\mathbf{M}_N \mathbf{V}^{(N)})^{\top} \mathbf{M}_N \mathbf{V}^{(N)} = \mathbf{V}^{(1)\top} \mathbf{V}^{(1)} \otimes \dots \otimes \mathbf{V}^{(N)\top} \mathbf{V}^{(N)}$ is diagonal. Hence, $\text{rank}_{\perp}((\mathbf{M}_1, \dots, \mathbf{M}_N) \cdot \mathcal{A}) \leq \text{rank}_{\perp}(\mathcal{A})$.

On the other hand, we have

$$\mathcal{A} = \left(\mathbf{M}_1^{\top}, \dots, \mathbf{M}_N^{\top} \right) \cdot [(\mathbf{M}_1, \dots, \mathbf{M}_N) \cdot \mathcal{A}]$$

and hence $\text{rank}_{\perp}(\mathcal{A}) \leq \text{rank}_{\perp}((\mathbf{M}_1, \dots, \mathbf{M}_N) \cdot \mathcal{A})$. Combining these two parts yields the result. □

In [22, (2.8)], an upper bound of $\text{rank}_\perp(\mathcal{A})$ is given as

$$\text{rank}_\perp(\mathcal{A}) \leq \min_{m=1, \dots, N} \prod_{n \neq m} I_n.$$

We refine this result in terms of the multilinear rank. The n -rank of \mathcal{A} , denoted by $\text{rank}_n(\mathcal{A})$, is the dimension of the vector space spanned by all n -mode fibers, i.e., $\text{rank}_n(\mathcal{A}) = \text{rank}(\mathbf{A}_{(n)})$. The N -tuple $(\text{rank}_1(\mathcal{A}), \dots, \text{rank}_N(\mathcal{A}))$ is called the multilinear rank of \mathcal{A} .

Proposition 3.5 *Let $\mathcal{A} \in \mathbb{R}^{I_1 \times \dots \times I_N}$. Then*

$$\text{rank}_\perp(\mathcal{A}) \leq \min_{m=1, \dots, N} \prod_{n \neq m} \text{rank}_n(\mathcal{A}).$$

Proof Suppose \mathcal{A} has the following higher-order singular value decomposition (HOSVD) [8]:

$$\mathcal{A} = (\mathbf{U}_1, \dots, \mathbf{U}_N) \cdot \mathcal{S},$$

where $\mathbf{U}_n \in \mathbb{R}^{I_n \times I_n}$ is orthogonal and $s_{i_1 i_2 \dots i_N} = 0$ if there exists at least one $n \in \{1, \dots, N\}$ such that $i_n > \text{rank}_n(\mathcal{A})$. It follows from Lemma 3.4 that $\text{rank}_\perp(\mathcal{A}) = \text{rank}_\perp(\mathcal{S})$. Note that

$$\mathcal{S} = \sum_{i_k, k \neq m} \mathbf{e}_{i_1} \otimes \dots \otimes \mathbf{e}_{i_{m-1}} \otimes \mathcal{S}(i_1, \dots, i_{m-1}, :, i_{m+1}, \dots, i_N) \otimes \mathbf{e}_{i_{m+1}} \otimes \dots \otimes \mathbf{e}_{i_N},$$

is an orthogonal decomposition, where $\mathbf{e}_{i_k} \in \mathbb{R}^{I_k}$ is the standard basis vector and $\mathcal{S}(i_1, \dots, i_{m-1}, :, i_{m+1}, \dots, i_N)$ is a mode- m fiber. Hence $\text{rank}_\perp(\mathcal{S})$ is less than the number of non-zero mode- m fibers, which is at most $\prod_{n \neq m} \text{rank}_n(\mathcal{A})$. \square

In contrast to Proposition 2.1, we have the following proposition for orthogonal rank.

Proposition 3.6 *For any $R > 0$, the set $\{\mathcal{A} \in \mathbb{R}^{I_1 \times \dots \times I_N} : \text{rank}_\perp(\mathcal{A}) \leq R\}$ is closed in the normed space $\mathbb{R}^{I_1 \times \dots \times I_N}$. That is, the function $\text{rank}_\perp(\mathcal{A})$ is lower semicontinuous.*

Proof Suppose $\lim_{m \rightarrow \infty} \mathcal{A}_m = \mathcal{A}$, where $\text{rank}_\perp(\mathcal{A}_m) \leq R$. It suffices to prove that $\text{rank}_\perp(\mathcal{A}) \leq R$. Since $\text{rank}_\perp(\mathcal{A}_m) \leq R$, we can write

$$\mathcal{A}_m = \sum_{r=1}^R \sigma_{m,r} \mathcal{U}_{m,r}, \quad \mathcal{U}_{m,r} = \mathbf{u}_{m,r}^{(1)} \otimes \dots \otimes \mathbf{u}_{m,r}^{(N)},$$

where $\langle \mathcal{U}_{m,s}, \mathcal{U}_{m,t} \rangle = 0$ for all $s \neq t$ and $\|\mathbf{u}_{m,r}^{(n)}\| = 1$ for all $n = 1, \dots, N$ and $r = 1, \dots, R$. (If $\text{rank}_\perp(\mathcal{A}_m) < R$, we just need to set $\sigma_{m,r} = 0$ for $r = \text{rank}_\perp(\mathcal{A}_m) + 1, \dots, R$.) Then

$$\sum_{r=1}^R \sigma_{m,r}^2 = \|\mathcal{A}_m\|^2.$$

Since $\lim_{m \rightarrow \infty} \|\mathcal{A}_m\| = \|\mathcal{A}\|$, $\sigma_{m,r}$ are uniformly bounded. Thus we can find a subsequence with convergence $\lim_{k \rightarrow \infty} \sigma_{m_k,r} = \sigma_r$, $\lim_{k \rightarrow \infty} \mathbf{u}_{m_k,r}^{(n)} = \mathbf{u}_r^{(n)}$ for all r and n . Note that $\lim_{k \rightarrow \infty} \mathcal{A}_{m_k} = \mathcal{A}$, i.e.,

$$\mathcal{A} = \sum_{r=1}^R \sigma_r \mathbf{u}_r^{(1)} \otimes \dots \otimes \mathbf{u}_r^{(N)}. \tag{8}$$

Moreover,

$$\left\langle \mathbf{u}_s^{(1)} \otimes \cdots \otimes \mathbf{u}_s^{(N)}, \mathbf{u}_t^{(1)} \otimes \cdots \otimes \mathbf{u}_t^{(N)} \right\rangle = \lim_{k \rightarrow \infty} \langle \mathcal{U}_{m_k, s}, \mathcal{U}_{m_k, t} \rangle = 0 \text{ for all } s \neq t.$$

Then (8) is an orthogonal decomposition and hence $\text{rank}_\perp(\mathcal{A}) \leq R$. □

Given $R > 0$, finding the *best orthogonal rank- R approximation* of \mathcal{A} is

$$\min_{\text{rank}_\perp(\mathcal{B}) \leq R} \|\mathcal{A} - \mathcal{B}\|. \tag{9}$$

By Proposition 3.6, we know that the solution of (9) always exists.

4 Algorithms for Low Orthogonal Rank Approximation

Problem (9) can be formulated as

$$\begin{aligned} \min_{\mathbf{v} \in \mathbb{R}^P} \mathcal{F}(\mathbf{v}) &:= \frac{1}{2} \left\| \mathcal{A} - \sum_{r=1}^R \otimes_{n=1}^N \mathbf{v}_r^{(n)} \right\|^2 \\ \text{s.t. } \prod_{n=1}^N \langle \mathbf{v}_s^{(n)}, \mathbf{v}_t^{(n)} \rangle &= 0 \text{ for all } s \neq t, \end{aligned} \tag{10}$$

where $\mathbf{v} := [\mathbf{v}_1^{(1)\top} \cdots \mathbf{v}_R^{(1)\top} \cdots \mathbf{v}_1^{(N)\top} \cdots \mathbf{v}_R^{(N)\top}]^\top$ and $P = R \sum_{n=1}^N I_n$. Define $\mathcal{T}_r = \otimes_{n=1}^N \mathbf{v}_r^{(n)}$, $r = 1, \dots, R$. Then (10) can be rewritten as

$$\min_{\mathbf{v} \in \mathbb{R}^P} \mathcal{F}(\mathbf{v}) := \frac{1}{2} \left\| \mathcal{A} - \sum_{r=1}^R \mathcal{T}_r \right\|^2 \text{ s.t. } \langle \mathcal{T}_s, \mathcal{T}_t \rangle = 0 \text{ for all } s \neq t. \tag{11}$$

The main difficulty in solving (10) is how to handle the $2N$ -degree polynomial constraints. The augmented Lagrangian method is a powerful method for solving constrained problems. This method is suitable for (10), because the subproblems can be solved easily by gradient-based optimization methods, which will be shown later. Using the augmented Lagrangian method means that the orthogonality constraint is not met exactly in each step. Our target is to make $\angle(\mathcal{T}_s, \mathcal{T}_t)$ close to $\pi/2$. By (3), we have

$$|\langle \mathcal{T}_s, \mathcal{T}_t \rangle| = \|\mathcal{T}_s\| \|\mathcal{T}_t\| |\cos \angle(\mathcal{T}_s, \mathcal{T}_t)|.$$

To avoid the influence of the norms $\|\mathcal{T}_r\|$, an ideal strategy is to consider the following augmented Lagrangian function:

$$\begin{aligned} \tilde{\mathcal{L}}(\mathbf{v}, \boldsymbol{\lambda}; \mathbf{c}) &= \mathcal{F}(\mathbf{v}) + \frac{1}{2} \sum_{s=1}^R \sum_{t=1, t \neq s}^R \lambda_{st} \prod_{n=1}^N \left\langle \frac{\mathbf{v}_s^{(n)}}{\|\mathbf{v}_s^{(n)}\|}, \frac{\mathbf{v}_t^{(n)}}{\|\mathbf{v}_t^{(n)}\|} \right\rangle \\ &+ \frac{\mu}{4} \sum_{s=1}^R \sum_{t=1, t \neq s}^R \prod_{n=1}^N \left\langle \frac{\mathbf{v}_s^{(n)}}{\|\mathbf{v}_s^{(n)}\|}, \frac{\mathbf{v}_t^{(n)}}{\|\mathbf{v}_t^{(n)}\|} \right\rangle^2. \end{aligned} \tag{12}$$

However, this would make the subproblem rather difficult to solve. We can realize this idea by setting different penalty parameters for the augmented Lagrangian function:

$$\begin{aligned} \mathcal{L}(\mathbf{v}, \boldsymbol{\lambda}; \mathbf{c}) := & \mathcal{F}(\mathbf{v}) + \frac{1}{2} \sum_{s=1}^R \sum_{t=1, t \neq s}^R \lambda_{st} \prod_{n=1}^N \langle \mathbf{v}_s^{(n)}, \mathbf{v}_t^{(n)} \rangle \\ & + \frac{1}{4} \sum_{s=1}^R \sum_{t=1, t \neq s}^R c_{st} \prod_{n=1}^N \langle \mathbf{v}_s^{(n)}, \mathbf{v}_t^{(n)} \rangle^2, \end{aligned} \tag{13}$$

where $\lambda_{st} = \lambda_{ts}$ are Lagrange multipliers, $c_{st} = c_{ts} > 0$ are penalty parameters and $\boldsymbol{\lambda} = \{\lambda_{st}\}$, $\mathbf{c} = \{c_{st}\}$. How to set penalty parameters will be elaborated in Sect. 4.1. Using different penalty parameters can reflect different features of different constraints, and is very common for the augmented Lagrangian method; see [3, p. 124], [33, Chapter 10.4] and [7, (1.5)].

For each iteration of the augmented Lagrangian method, we need to solve the following problem

$$\min_{\mathbf{v} \in \mathbb{R}^P} \mathcal{L}(\mathbf{v}, \boldsymbol{\lambda}; \mathbf{c}) \tag{14}$$

with $\boldsymbol{\lambda}, \mathbf{c}$ given. If $\boldsymbol{\lambda} = \{0\}$, $\mathbf{c} = \{0\}$, (14) is just (6). Since (6) has no solution in general, the first issue that we need to make sure is whether (14) has a solution. We have the following proposition.

Proposition 4.1 *If $c_{st} > 0$ for all $s \neq t$, then (14) always has a solution.*

Proof For convenience, define $\mathcal{E}(\mathbf{v}) = \mathcal{L}(\mathbf{v}, \boldsymbol{\lambda}; \mathbf{c})$. It follows from (11) and (13) that

$$\mathcal{E}(\mathbf{v}) = \frac{1}{2} \left\| \mathcal{A} - \sum_{r=1}^R \mathcal{T}_r \right\|^2 + \frac{1}{4} \sum_{s=1}^R \sum_{t=1, t \neq s}^R c_{st} \left(\langle \mathcal{T}_s, \mathcal{T}_t \rangle + \frac{\lambda_{st}}{c_{st}} \right)^2 - \frac{1}{4} \sum_{s=1}^R \sum_{t=1, t \neq s}^R \frac{\lambda_{st}^2}{c_{st}}.$$

Note that

$$\otimes_{n=1}^N \mathbf{v}_r^{(n)} = \otimes_{n=1}^N b_n \mathbf{v}_r^{(n)} \quad \text{when} \quad \prod_{n=1}^N b_n = 1. \tag{15}$$

We can scale each $\mathbf{v}_r^{(n)}$ such that $\|\mathbf{v}_r^{(n)}\| = \|\mathcal{T}_r\|^{1/N}$, $n = 1, \dots, N$. Define the following set

$$W = \{ \mathbf{v} \in \mathbb{R}^P : \|\mathbf{v}_r^{(m)}\| = \|\mathbf{v}_r^{(n)}\|, 1 \leq m, n \leq N, 1 \leq r \leq R \}.$$

The continuity of $\|\cdot\|$ implies that W is closed. We have

$$\{ \mathcal{E}(\mathbf{v}) : \mathbf{v} \in \mathbb{R}^P \} = \{ \mathcal{E}(\mathbf{v}) : \mathbf{v} \in W \}.$$

Hence, it suffices to show that (14) has a solution on W .

Define $\alpha = \frac{1}{4} \sum_{s=1}^R \sum_{t=1, t \neq s}^R \frac{\lambda_{st}^2}{c_{st}}$, $\beta = \min\{c_{st}\}$, $\gamma = \sum_{s=1}^R \sum_{t=1, t \neq s}^R \frac{|\lambda_{st}|}{c_{st}}$. For any $\xi \geq \inf \mathcal{E} \geq 0$, if $\mathcal{E} \leq \xi$, then $\|\mathcal{A} - \sum_{r=1}^R \mathcal{T}_r\| \leq \sqrt{2(\xi + \alpha)}$ and

$$\begin{aligned} & \sum_{s=1}^R \sum_{t=1, t \neq s}^R |\langle \mathcal{T}_s, \mathcal{T}_t \rangle| - \gamma \leq \sum_{s=1}^R \sum_{t=1, t \neq s}^R \left| \langle \mathcal{T}_s, \mathcal{T}_t \rangle + \frac{\lambda_{st}}{c_{st}} \right| \\ & \leq \sqrt{R(R-1) \sum_{s=1}^R \sum_{t=1, t \neq s}^R \left(\langle \mathcal{T}_s, \mathcal{T}_t \rangle + \frac{\lambda_{st}}{c_{st}} \right)^2} \leq \sqrt{\frac{4R(R-1)(\xi + \alpha)}{\beta}} \\ & \implies \sum_{s=1}^R \sum_{t=1, t \neq s}^R |\langle \mathcal{T}_s, \mathcal{T}_t \rangle| \leq \gamma + \sqrt{\frac{4R(R-1)(\xi + \alpha)}{\beta}}. \end{aligned}$$

Hence $\|\sum_{r=1}^R \mathcal{T}_r\| \leq \|\mathcal{A} - \sum_{r=1}^R \mathcal{T}_r\| + \|\mathcal{A}\| \leq \sqrt{2(\xi + \alpha)} + \|\mathcal{A}\|$. For any $\mathbf{v} \in W$, it follows that

$$\begin{aligned} & (\sqrt{2(\xi + \alpha)} + \|\mathcal{A}\|)^2 \geq \left\| \sum_{r=1}^R \mathcal{T}_r \right\|^2 = \sum_{r=1}^R \|\mathcal{T}_r\|^2 + \sum_{s=1}^R \sum_{t=1, t \neq s}^R \langle \mathcal{T}_s, \mathcal{T}_t \rangle \\ & \geq \sum_{r=1}^R \|\mathcal{T}_r\|^2 - \sum_{s=1}^R \sum_{t=1, t \neq s}^R |\langle \mathcal{T}_s, \mathcal{T}_t \rangle| \geq \sum_{r=1}^R \|\mathcal{T}_r\|^2 - \sqrt{\frac{4R(R-1)(\xi + \alpha)}{\beta}} - \gamma \\ & \implies \|\mathbf{v}_r^n\|^2 = \|\mathcal{T}_r\|^2 \leq \left((\sqrt{2(\xi + \alpha)} + \|\mathcal{A}\|)^2 + \sqrt{\frac{4R(R-1)(\xi + \alpha)}{\beta}} + \gamma \right)^{1/N}. \end{aligned}$$

That is, the level set $\{\mathbf{v} \in W : \mathcal{E}(\mathbf{v}) \leq \xi, \xi \geq \inf \mathcal{E}\}$ is bounded. Combining with the fact that $\mathcal{E}(\mathbf{v})$ is continuous and W is closed, it follows from [28, Theorem 1.9] that \mathcal{E} can attain its minimum on W . □

We will employ gradient-based optimization methods to solve each subproblem. Gradient-based optimization methods have been used in fitting CP decompositions; see [1, 12]. To use such methods, we need to compute the gradient of the objective function. Define

$$\mathbf{V}^{(-n)} := \mathbf{V}^{(N)} \odot \dots \odot \mathbf{V}^{(n+1)} \odot \mathbf{V}^{(n-1)} \odot \dots \odot \mathbf{V}^{(1)}, \tag{16}$$

where “ \odot ” is the Khatri-Rao product. With the relationship introduced in [19, Section 2.6], we have

$$\mathbf{\Gamma}^{(n)} := \mathbf{V}^{(-n)\top} \mathbf{V}^{(-n)} = \left(\otimes_{m=1}^{n-1} \left(\mathbf{V}^{(m)\top} \mathbf{V}^{(m)} \right) \right) \otimes \left(\otimes_{m=n+1}^N \left(\mathbf{V}^{(m)\top} \mathbf{V}^{(m)} \right) \right). \tag{17}$$

The gradient of the first term of \mathcal{L} , i.e., $\mathcal{F}(\mathbf{v})$ in (10), can be found in [1, 12]. Here we provide a calculation based on unfolding matrices as follows.

Lemma 4.2 *The gradient of $\mathcal{F}(\mathbf{v})$ in (10) is given by*

$$\frac{\partial \mathcal{F}}{\partial \mathbf{V}^{(n)}} = \left[\frac{\partial \mathcal{F}}{\partial \mathbf{v}_1^{(n)}} \dots \frac{\partial \mathcal{F}}{\partial \mathbf{v}_R^{(n)}} \right] = -\mathbf{A}^{(n)} \mathbf{V}^{(-n)} + \mathbf{V}^{(n)} \mathbf{\Gamma}^{(n)}$$

for $n = 1, \dots, N$, where $\mathbf{V}^{(-n)}$ and $\mathbf{\Gamma}^{(n)}$ are defined in (16) and (17), respectively.

Proof Let $\mathcal{B} = \sum_{r=1}^R \otimes_{n=1}^N \mathbf{v}_r^{(n)}$. Then we have $\mathbf{B}_{(n)} = \mathbf{V}^{(n)} \mathbf{V}^{(-n)\top}$ (see [19, Sect. 3]) and

$$\mathcal{F}(\mathbf{v}) = \frac{1}{2} \|\mathcal{A} - \mathcal{B}\|^2 = \frac{1}{2} \left\| -\mathbf{A}_{(n)} + \mathbf{V}^{(n)} \mathbf{V}^{(-n)\top} \right\|^2.$$

Then,

$$\frac{\partial \mathcal{F}}{\partial \mathbf{V}^{(n)}} = \left(-\mathbf{A}_{(n)} + \mathbf{V}^{(n)} \mathbf{V}^{(-n)\top} \right) \mathbf{V}^{(-n)} = -\mathbf{A}_{(n)} \mathbf{V}^{(-n)} + \mathbf{V}^{(n)} \mathbf{\Gamma}^{(n)}.$$

□

Denote the sum of the last two terms of \mathcal{L} by \mathcal{G} , i.e.,

$$\mathcal{G} = \frac{1}{2} \sum_{s=1}^R \sum_{t=1, t \neq s}^R \lambda_{st} \prod_{n=1}^N \left\langle \mathbf{v}_s^{(n)}, \mathbf{v}_t^{(n)} \right\rangle + \frac{1}{4} \sum_{s=1}^R \sum_{t=1, t \neq s}^R c_{st} \prod_{n=1}^N \left\langle \mathbf{v}_s^{(n)}, \mathbf{v}_t^{(n)} \right\rangle^2,$$

and define $\gamma_{sr}^{(n)} = \prod_{m=1, m \neq n}^N \left\langle \mathbf{v}_s^{(m)}, \mathbf{v}_r^{(m)} \right\rangle$. Direct calculation gives that

$$\frac{\partial \mathcal{G}}{\partial \mathbf{v}_r^{(n)}} = \sum_{s=1, s \neq r}^R \left(\lambda_{sr} \gamma_{sr}^{(n)} + c_{sr} \gamma_{sr}^{(n)2} \left\langle \mathbf{v}_s^{(n)}, \mathbf{v}_r^{(n)} \right\rangle \right) \mathbf{v}_s^{(n)}. \tag{18}$$

Note that $\gamma_{st}^{(n)} = \mathbf{\Gamma}^{(n)}(s, t)$, where $\mathbf{\Gamma}^{(n)}$ is defined in (17). Define matrices $\mathbf{\Lambda}, \mathbf{C} \in \mathbb{R}^{R \times R}$ by

$$\mathbf{\Lambda}(i, j) = \begin{cases} \lambda_{ij}, & \text{if } i \neq j \\ 0, & \text{otherwise,} \end{cases} \quad \mathbf{C}(i, j) = \begin{cases} c_{ij}, & \text{if } i \neq j \\ 0, & \text{otherwise.} \end{cases} \tag{19}$$

Combining Lemma 4.2 and (18), we can write the gradient of \mathcal{L} in matrix form, as the following corollary shows.

Corollary 4.3 *The gradient of the objective function \mathcal{L} in (13) is given by*

$$\frac{\partial \mathcal{L}}{\partial \mathbf{V}^{(n)}} = -\mathbf{A}_{(n)} \mathbf{V}^{(-n)} + \mathbf{V}^{(n)} \left(\mathbf{\Gamma}^{(n)} + \mathbf{\Gamma}^{(n)} \circledast \mathbf{\Lambda} + \mathbf{\Gamma}^{(n)} \circledast \mathbf{\Gamma}^{(n)} \circledast \mathbf{V}^{(n)\top} \mathbf{V}^{(n)} \circledast \mathbf{C} \right)$$

for $n = 1, \dots, N$, where related matrices are defined in (16), (17) and (19).

4.1 Algorithm: OD-ALM

Suppose we have obtained the solution $\mathbf{v}_{[k]}$ for the k th iteration. We introduce how to solve $\mathbf{v}_{[k+1]}$ for the $(k + 1)$ st iteration.

We use $\mathbf{v}_{[k]}$ as the initialization of the $(k + 1)$ st iteration. By (15), we scale the initialization such that $\|\mathbf{v}_{r,[k]}^{(m)}\| = \left(\prod_{n=1}^N \|\mathbf{v}_{r,[k]}^{(n)}\| \right)^{1/N}$, $m = 1, \dots, N$. This scaling can avoid the situation that some $\|\mathbf{v}_{r,[k]}^{(n_1)}\|$ is too big and some $\|\mathbf{v}_{r,[k]}^{(n_2)}\|$ is too small, where $1 \leq n_1, n_2 \leq N$. The idea of (12) can be realized by setting different penalty parameters for (13):

$$c_{st,[k]} = \frac{\mu_{[k]}}{\prod_{n=1}^N \|\mathbf{v}_{s,[k]}^{(n)}\|^2 \prod_{n=1}^N \|\mathbf{v}_{t,[k]}^{(n)}\|^2}, \tag{20}$$

where $\mu_{[k]} > 0$. In the matrix form (19), the non-diagonal entries of $\mathbf{C}_{[k]}$ are the same as those of $\mu_{[k]} \mathbf{h}_{[k]}^\top \mathbf{h}_{[k]}$, where

$$\mathbf{h}_{[k]} = \left[\frac{1}{\prod_{n=1}^N \|\mathbf{v}_{1,[k]}^{(n)}\|^2} \cdots \frac{1}{\prod_{n=1}^N \|\mathbf{v}_{R,[k]}^{(n)}\|^2} \right] \in \mathbb{R}^{1 \times R}.$$

Then $\mathbf{v}_{[k+1]}$ can be obtained by solving $\min_{\mathbf{v} \in \mathbb{R}^P} \mathcal{L}(\mathbf{v}, \boldsymbol{\lambda}_{[k]}; \mathbf{c}_{[k]})$. At last, the Lagrange multiplier $\lambda_{st,[k+1]}$ is updated by $\lambda_{st,[k+1]} = \lambda_{st,[k]} + c_{st,[k]} \prod_{n=1}^N \langle \mathbf{v}_{s,[k+1]}^{(n)}, \mathbf{v}_{t,[k+1]}^{(n)} \rangle$, whose matrix form is

$$\mathbf{\Lambda}_{[k+1]} = \mathbf{\Lambda}_{[k]} + \mathbf{C}_{[k]} \circledast \left(\otimes_{n=1}^N \mathbf{V}_{[k+1]}^{(n)\top} \mathbf{V}_{[k+1]}^{(n)} \right). \tag{21}$$

Now we introduce how to develop a systematic scheme for the augmented Lagrangian method. The standard procedure of the augmented Lagrangian method tells us that we need to increase the penalty parameters gradually to a sufficiently large value. This procedure is rather important for (14), because \mathcal{L} is nonconvex. The later subproblems corresponding to larger penalty parameters can be solved relatively efficiently by warm starting point from the previous solutions. By (20), we need to set $\mu_{[k+1]}$ sufficiently large such that $c_{st,[k+1]} > c_{st,[k]}$ for all $s \neq t$. We set

$$\mu_{[k+1]} = \max \left\{ \beta, \max_{s \neq t} \left\{ \prod_{n=1}^N \frac{\|\mathbf{v}_{s,[k+1]}^{(n)}\|^2 \|\mathbf{v}_{t,[k+1]}^{(n)}\|^2}{\|\mathbf{v}_{s,[k]}^{(n)}\|^2 \|\mathbf{v}_{t,[k]}^{(n)}\|^2} \right\} \right\} \mu_{[k]}, \tag{22}$$

where $\beta > 1$. Usually, we can simply set a sufficiently large β , for instance $\beta = 10$, and the condition $c_{st,[k+1]} > c_{st,[k]}$ will be satisfied naturally. The whole procedure of the augmented Lagrangian method is presented in Algorithm 1.

Algorithm 1: Orthogonal Decomposition by Augmented Lagrangian Method (OD-ALM)

```

Input: Tensor  $\mathcal{A}$ , number of components  $R$ , initialization  $\mathbf{v}_{[0]}$ ;  $\mathbf{\Lambda}_{[0]} = \mathbf{0}, \mu_{[0]} = 1; k = 0$ 
Output: Approximate solution  $\mathbf{v}_{[k]}$  of the orthogonal rank- $R$  approximation to  $\mathcal{A}$ 
1 repeat
2   for  $r = 1, \dots, R$  do
3      $\delta_r \leftarrow \prod_{n=1}^N \|\mathbf{v}_{r,[k]}^{(n)}\|$  ▷ Compute the norm of  $\otimes_{n=1}^N \mathbf{v}_{r,[k]}^{(n)}$ 
4   end
5   for  $r = 1, \dots, R$  do
6     for  $n = 1, \dots, N$  do
7        $\mathbf{v}_{r,[k]}^{(n)} \leftarrow \frac{\delta_r^{1/N}}{\|\mathbf{v}_{r,[k]}^{(n)}\|} \mathbf{v}_{r,[k]}^{(n)}$  ▷ scale the initialization
8     end
9   end
10   $\mathbf{h} \leftarrow [1/\delta_1^2 \ \dots \ 1/\delta_R^2]$ 
11   $\mathbf{C}_{[k]} \leftarrow \mu_{[k]} \mathbf{h}^\top \mathbf{h}$ 
12   $\mathbf{C}_{[k]}(i, i) \leftarrow 0 \ \forall i = 1, \dots, R$ 
13   $\mathbf{v}_{[k+1]} \leftarrow \arg \min \mathcal{L}(\mathbf{v}, \boldsymbol{\lambda}_{[k]}; \mathbf{c}_{[k]})$  by gradient-based optimization methods with starting point  $\mathbf{v}_{[k]}$ , where the gradient is computed by Corollary 4.3
14  Update  $\mathbf{\Lambda}_{[k+1]}$  by (21)
15  Update  $\mu_{[k+1]}$  by (22)
16   $k \leftarrow k + 1$ 
17 until termination criteria met

```

The convergence analysis of augmented Lagrangian methods can be found in many textbooks. See [3, 27, 33] for reference. Here we extend [33, Theorem 10.4.2], which is useful for designing the termination criteria.

Proposition 4.4 For Algorithm 1, we have

$$\lim_{k \rightarrow \infty} \prod_{n=1}^N \left\langle \frac{\mathbf{v}_{s,[k+1]}^{(n)}}{\|\mathbf{v}_{s,[k]}^{(n)}\|}, \frac{\mathbf{v}_{t,[k+1]}^{(n)}}{\|\mathbf{v}_{t,[k]}^{(n)}\|} \right\rangle = 0 \text{ for all } 1 \leq s \neq t \leq R.$$

Proof We have

$$\begin{aligned} \sum_{s \neq t} \frac{\lambda_{st,[k+1]}^2}{c_{st,[k+1]}} &\leq \sum_{s \neq t} \frac{\lambda_{st,[k+1]}^2}{c_{st,[k]}} = \sum_{s \neq t} \frac{\left(\lambda_{st,[k]} + c_{st,[k]} \prod_n \left\langle \frac{\mathbf{v}_{s,[k+1]}^{(n)}}{\|\mathbf{v}_{s,[k]}^{(n)}\|}, \frac{\mathbf{v}_{t,[k+1]}^{(n)}}{\|\mathbf{v}_{t,[k]}^{(n)}\|} \right\rangle \right)^2}{c_{st,[k]}} \\ &= \sum_{s \neq t} \frac{\lambda_{st,[k]}^2}{c_{st,[k]}} + 4 \left(\mathcal{L}(\mathbf{v}_{[k+1]}, \boldsymbol{\lambda}_{[k]}; \mathbf{c}_{[k]}) - \mathcal{F}(\mathbf{v}_{[k+1]}) \right) \\ &\leq \sum_{s \neq t} \frac{\lambda_{st,[k]}^2}{c_{st,[k]}} + 4 \mathcal{L}(\mathbf{v}_{[k+1]}, \boldsymbol{\lambda}_{[k]}; \mathbf{c}_{[k]}). \end{aligned}$$

For any feasible point $\bar{\mathbf{v}}$ of (10), by noting that $\mathcal{L}(\mathbf{v}_{[k+1]}, \boldsymbol{\lambda}_{[k]}; \mathbf{c}_{[k]}) \leq \mathcal{L}(\bar{\mathbf{v}}, \boldsymbol{\lambda}_{[k]}; \mathbf{c}_{[k]}) = \mathcal{F}(\bar{\mathbf{v}})$, we have

$$\sum_{s \neq t} \frac{\lambda_{st,[k+1]}^2}{c_{st,[k+1]}} \leq \sum_{s \neq t} \frac{\lambda_{st,[k]}^2}{c_{st,[k]}} + 4 \mathcal{L}(\mathbf{v}_{[k+1]}, \boldsymbol{\lambda}_{[k]}; \mathbf{c}_{[k]}) \leq \sum_{s \neq t} \frac{\lambda_{st,[k]}^2}{c_{st,[k]}} + 4 \mathcal{F}(\bar{\mathbf{v}}).$$

This suggests that there exists $\delta > 0$ such that $\sum_{s \neq t} \frac{\lambda_{st,[k]}^2}{c_{st,[k]}} \leq \delta k$. Define

$$d_{st,[k]} := \lambda_{st,[k]} \prod_n \|\mathbf{v}_{s,[k]}^{(n)}\| \prod_n \|\mathbf{v}_{t,[k]}^{(n)}\|.$$

It follows from (20) that $\sum_{s \neq t} \frac{d_{st,[k]}^2}{\mu_{[k]}} = \sum_{s \neq t} \frac{\lambda_{st,[k]}^2}{c_{st,[k]}} \leq \delta k$. By the algorithm, $\mu_{[k]} \geq \beta^k$, where $\beta > 1$. Hence, $\frac{d_{st,[k]}}{\mu_{[k]}} = o(1)$.

For any feasible point $\bar{\mathbf{v}}$ of (10), we have

$$\begin{aligned} \mathcal{F}(\bar{\mathbf{v}}) &= \mathcal{L}(\bar{\mathbf{v}}, \boldsymbol{\lambda}_{[k]}; \mathbf{c}_{[k]}) \geq \mathcal{L}(\mathbf{v}_{[k+1]}, \boldsymbol{\lambda}_{[k]}; \mathbf{c}_{[k]}) \\ &= \mathcal{F}(\mathbf{v}_{[k+1]}) + \frac{1}{2} \sum_{s \neq t} d_{st,[k]} \prod_{n=1}^N \left\langle \frac{\mathbf{v}_{s,[k+1]}^{(n)}}{\|\mathbf{v}_{s,[k]}^{(n)}\|}, \frac{\mathbf{v}_{t,[k+1]}^{(n)}}{\|\mathbf{v}_{t,[k]}^{(n)}\|} \right\rangle \\ &\quad + \frac{1}{4} \sum_{s \neq t} \mu_{[k]} \prod_{n=1}^N \left\langle \frac{\mathbf{v}_{s,[k+1]}^{(n)}}{\|\mathbf{v}_{s,[k]}^{(n)}\|}, \frac{\mathbf{v}_{t,[k+1]}^{(n)}}{\|\mathbf{v}_{t,[k]}^{(n)}\|} \right\rangle^2 \\ &= \mathcal{F}(\mathbf{v}_{[k+1]}) + \frac{1}{4} \sum_{s \neq t} \mu_{[k]} \left[\left(\prod_{n=1}^N \left\langle \frac{\mathbf{v}_{s,[k+1]}^{(n)}}{\|\mathbf{v}_{s,[k]}^{(n)}\|}, \frac{\mathbf{v}_{t,[k+1]}^{(n)}}{\|\mathbf{v}_{t,[k]}^{(n)}\|} \right\rangle + \frac{d_{st,[k]}}{\mu_{[k]}} \right)^2 - \left(\frac{d_{st,[k]}}{\mu_{[k]}} \right)^2 \right] \\ &\geq \frac{1}{4} \sum_{s \neq t} \mu_{[k]} \left[\left(\prod_{n=1}^N \left\langle \frac{\mathbf{v}_{s,[k+1]}^{(n)}}{\|\mathbf{v}_{s,[k]}^{(n)}\|}, \frac{\mathbf{v}_{t,[k+1]}^{(n)}}{\|\mathbf{v}_{t,[k]}^{(n)}\|} \right\rangle + o(1) \right)^2 - o(1) \right]. \end{aligned}$$

By noting that $\lim_{k \rightarrow \infty} \mu_{[k]} = \infty$ and $\mathcal{F}(\bar{\mathbf{v}})$ is bounded, we obtain the result. \square

Corollary 4.5 For Algorithm 1, suppose $\prod_{n=1}^N \frac{\|\mathbf{v}_{r,[k]}^{(n)}\|}{\|\mathbf{v}_{r,[k+1]}^{(n)}\|}$ is bounded for all r and k . Then we have

$$\lim_{k \rightarrow \infty} \prod_{n=1}^N \left\langle \frac{\mathbf{v}_{s,[k]}^{(n)}}{\|\mathbf{v}_{s,[k]}^{(n)}\|}, \frac{\mathbf{v}_{t,[k]}^{(n)}}{\|\mathbf{v}_{t,[k]}^{(n)}\|} \right\rangle = 0 \text{ for all } 1 \leq s \neq t \leq R.$$

4.2 Orthogonalization of Rank-One Tensors

OD-ALM can only obtain an approximate solution of (10). We need to develop an orthogonalization process to make the orthogonality constraint exact for the final result.

Suppose we have obtained a decomposition by OD-ALM:

$$\mathcal{A} \approx \sum_{r=1}^R \otimes_{n=1}^N \mathbf{v}_r^{(n)}.$$

First, we normalize each $\mathbf{v}_r^{(n)}$ to $\mathbf{u}_r^{(n)}$, i.e., $\mathbf{u}_r^{(n)} = \mathbf{v}_r^{(n)} / \|\mathbf{v}_r^{(n)}\|$. Assume that we have orthogonalized the first $\ell - 1$ rank-one components:

$$\left\langle \otimes_{n=1}^N \mathbf{u}_s^{(n)}, \otimes_{n=1}^N \mathbf{u}_t^{(n)} \right\rangle = 0, \quad 1 \leq s \neq t \leq \ell - 1.$$

We start to handle the ℓ th rank-one component. Define

$$\widehat{\mathbf{U}}^{(n)} := \left[\mathbf{u}_1^{(n)} \ \dots \ \mathbf{u}_{\ell-1}^{(n)} \right], \quad n = 1, \dots, N.$$

Compute the absolute value of the inner product $\left| \left\langle \mathbf{u}_r^{(n)}, \mathbf{u}_\ell^{(n)} \right\rangle \right|$ for $n = 1, \dots, N$ and $r = 1, \dots, \ell - 1$, and stack the results as a matrix:

$$\mathbf{P} = \left[\begin{array}{c} \left[\mathbf{u}_\ell^{(1)\top} \widehat{\mathbf{U}}^{(1)} \right] \\ \vdots \\ \left[\mathbf{u}_\ell^{(N)\top} \widehat{\mathbf{U}}^{(N)} \right] \end{array} \right] \in \mathbb{R}^{N \times (\ell-1)},$$

where $|\cdot|$ denotes the entrywise absolute value. Let $\mathbf{P}(m_r, r) = \min\{\mathbf{P}(1, r), \dots, \mathbf{P}(N, r)\}$. Then for each $r \in \{1, \dots, \ell - 1\}$, $\{\mathbf{u}_r^{(m_r)}, \mathbf{u}_\ell^{(m_r)}\}$ is a pair of vectors that is the closest to orthogonality among all pairs $\{\mathbf{u}_r^{(n)}, \mathbf{u}_\ell^{(n)}\}, n = 1, \dots, N$. Suppose $\{r : m_r = n\} = \{r_1, \dots, r_{\rho(n)}\}$. For each $n \in \{1, \dots, N\}$, We will modify $\mathbf{u}_\ell^{(n)}$ to $\mathbf{u}_\ell^{(n)} - \sum_{j=1}^{\rho(n)} x_j \mathbf{u}_{r_j}^{(n)}$ such that

$$\left\langle \mathbf{u}_\ell^{(n)} - \sum_{j=1}^{\rho(n)} x_j \mathbf{u}_{r_j}^{(n)}, \mathbf{u}_s^{(n)} \right\rangle = 0, \quad s = r_1, \dots, r_{\rho(n)},$$

whose matrix form is

$$\left[\mathbf{u}_{r_1}^{(n)} \ \dots \ \mathbf{u}_{r_{\rho(n)}}^{(n)} \right]^\top \left[\mathbf{u}_{r_1}^{(n)} \ \dots \ \mathbf{u}_{r_{\rho(n)}}^{(n)} \right] \begin{bmatrix} x_1 \\ \vdots \\ x_{\rho(n)} \end{bmatrix} = \left[\mathbf{u}_{r_1}^{(n)} \ \dots \ \mathbf{u}_{r_{\rho(n)}}^{(n)} \right]^\top \mathbf{u}_\ell^{(n)}.$$

We present the whole process of the orthogonalization in Algorithm 2. This process can also be used for generating general orthonormal lists of rank-one tensors.

Algorithm 2: Orthogonalization of rank-one tensors

```

Input: A list of rank-one tensors  $\{\mathbf{v}_r^{(n)}\}_{n,r}$ 
Output: An orthonormal list rank-one tensors  $\{\mathbf{u}_r^{(n)}\}_{n,r}$ 
1 for  $r = 1, \dots, R$  do
2   for  $n = 1, \dots, N$  do
3      $\eta \leftarrow \|\mathbf{v}_r^{(n)}\|$ 
4      $\mathbf{u}_r^{(n)} \leftarrow \mathbf{v}_r^{(n)} / \eta$ 
5   end
6 end
7 for  $\ell = 2, \dots, R$  do
8   for  $n = 1, \dots, N$  do
9      $\mathbf{U} \leftarrow [\mathbf{u}_1^{(n)} \dots \mathbf{u}_{\ell-1}^{(n)}]$ 
10     $\mathbf{P}(n, :) \leftarrow |\mathbf{u}_\ell^{(n)\top} \mathbf{U}|$ 
11  end
12  for  $r = 1, \dots, \ell - 1$  do
13    Find  $\mathbf{P}(m_r, r) = \min\{\mathbf{P}(1, r), \dots, \mathbf{P}(N, r)\}$ 
14  end
15  for  $n = 1, \dots, N$  do
16     $\{r_1, \dots, r_{\rho(n)}\} \leftarrow$  all indices satisfying  $m_{r_j} = n, j = 1, \dots, \rho(n)$ 
17    if  $\rho(n) = 0$  then
18       $\mathbf{u}_\ell^{(n)} \leftarrow \mathbf{u}_\ell^{(n)}$ 
19    else
20       $\mathbf{B} \leftarrow [\mathbf{u}_{r_1}^{(n)} \dots \mathbf{u}_{r_{\rho(n)}}^{(n)}]$ 
21      Solve  $\mathbf{B}^\top \mathbf{B} \mathbf{x} = \mathbf{B}^\top \mathbf{u}_\ell^{(n)}$  for  $\mathbf{x}$ 
22       $\mathbf{u}_\ell^{(n)} \leftarrow \mathbf{u}_\ell^{(n)} - \mathbf{B} \mathbf{x}$ 
23       $\eta \leftarrow \|\mathbf{u}_\ell^{(n)}\|$ 
24       $\mathbf{u}_r^{(n)} \leftarrow \mathbf{u}_r^{(n)} / \eta$ 
25    end
26  end
27 end

```

The final orthogonal rank- R approximation is the orthogonal projection of \mathcal{A} onto the space spanned by the orthonormal list $\{\otimes_{n=1}^N \mathbf{u}_1^{(n)}, \dots, \otimes_{n=1}^N \mathbf{u}_R^{(n)}\}$:

$$\sum_{r=1}^R \sigma_r \otimes_{n=1}^N \mathbf{u}_r^{(n)},$$

where the coefficient $\sigma_r = \langle \mathcal{A}, \otimes_{n=1}^N \mathbf{u}_r^{(n)} \rangle$.

5 Numerical Experiments

We will show the performance of OD-ALM combined with the orthogonalization process in this section. All experiments are performed on MATLAB R2016a with Tensor Toolbox, version 3.0 [2] on a laptop (Intel Core i5-6300HQ CPU @ 2.30GHz, 8.00G RAM). The test data include both synthetic and real-world tensors. The synthetic tensors are generated from

Table 1 The test tensors. The value R is the number of components for all methods

Tensor	Size	R	Note
\mathcal{A}_1	$20 \times 16 \times 10 \times 32$	5	random tensor
\mathcal{A}_2	$20 \times 16 \times 10 \times 32$	5	rank-5 tensor
\mathcal{A}_3	$20 \times 16 \times 10 \times 32$	5	$\mathcal{A}_3(i_1, i_2, i_3, i_4) = 1/(i_1 + i_2 + i_3 + i_4 - 3)$
\mathcal{A}_4	$20 \times 16 \times 10 \times 32$	5	orthogonal rank-5 tensor with Gaussian noise
\mathcal{A}_5	$95 \times 95 \times 156$	5	hyperspectral image – Samson
\mathcal{A}_6	$100 \times 100 \times 224$	5	hyperspectral image – Jasper Ridge
\mathcal{A}_7	$144 \times 176 \times 3 \times 300$	2	video data – Akiyo
\mathcal{A}_8	$144 \times 176 \times 3 \times 300$	2	video data – Hall Monitor

known ground truth and thus make the evaluation reliable. Choosing real-world tensors is to assess the approximation ability of orthogonal decompositions in practice.

The experiments are designed based on those of [5, 13]. The test tensors are shown in Table 1, where $\mathcal{A}_1, \dots, \mathcal{A}_4$ are synthetic tensors and $\mathcal{A}_5, \dots, \mathcal{A}_8$ are real-world tensors. The tensor \mathcal{A}_1 is a randomly generated tensor, \mathcal{A}_2 is a randomly generated rank-5 tensor, and \mathcal{A}_3 is a Hilbert tensor also used in [13]. For \mathcal{A}_4 , we generate an orthonormal list of rank-one tensors by Algorithm 2 and then use this list to generate an orthogonal rank-5 tensor \mathcal{B}_1 . The final tensor \mathcal{A}_4 is

$$\mathcal{A}_4 = \mathcal{B}_1 + \rho \mathcal{B}_2,$$

where \mathcal{B}_2 is a noise tensor with entries drawn from a standard normal distribution, and $\rho = 0.1 \|\mathcal{B}_1\| / \|\mathcal{B}_2\|$. The tensors $\mathcal{A}_5, \mathcal{A}_6$ are hyperspectral images³, and $\mathcal{A}_7, \mathcal{A}_8$ are video tensors⁴. We will factorize each tensor into R terms by different methods. Different R 's would result in different approximation errors and running time. We concern the approximation abilities of different methods. For simplicity, we fix R for each test tensor, prescribed in Table 1. Using other R 's would show the same comparison results.

Suppose \mathcal{B} is an approximation of \mathcal{A} obtained by any method. We use the relative error (RErr) to evaluate the result:

$$\text{RErr} = \frac{\|\mathcal{A} - \mathcal{B}\|}{\|\mathcal{A}\|}.$$

5.1 Implementation Details of OD-ALM

The initialization is crucial for OD-ALM. We adopt the result of the alternating least squares algorithm (CP-ALS) [4, 14, 19] for (6) as the initialization, because this result is just the numerical solution of (14) with Lagrange multipliers and penalty parameters equal to zero, which is relatively near to the solution of the first subproblem of OD-ALM. The CP-ALS is with the truncated HOSVD initialization, and terminates if the relative change in the function value is less than 10^{-6} or the number of iterations exceeds 500. As for (22), we set $\beta = 10$ for all tests.

Commonly used gradient-based optimization methods include the steepest descent method, the conjugate gradient method, the Broyden-Fletcher-Goldfarb-Shanno (BFGS)

³ The hyperspectral image data have been used in [36] and available at <https://rslab.ut.ac.ir/data>.

⁴ The video data are from the video trace library [29] and available at <http://trace.eas.asu.edu/yuv/>.

method and the limited-memory BFGS (L-BFGS) method. We have tried all these methods to solve the subproblems (14) and find that the L-BFGS method outperforms the other three ones. Hence, we use the L-BFGS method with $m = 20$ levels of memory in all tests. We stop the procedure of the L-BFGS method if the relative change between successive iterates is less than 10^{-8} , or the ℓ_2 norm of the gradient divided by the number of entries is less than ϵ_{inner} , which will be specified later. The maximum number of inner iterations is set to be 500. We adopt the Moré-Thuente line search [25] from MINPACK⁵. For all experiments, Moré-Thuente line search parameters used are as follows: 10^{-4} for the function value tolerance, 10^{-2} for the gradient norm tolerance, a starting search step length of 1 and a maximum of 20 iterations.

For the solution $\mathbf{v}_{[k]}$ of the k th subproblem, define

$$\theta_{[k]} := \max_{s \neq t} \min_n \left\langle \left\| \frac{\mathbf{v}_{s,[k]}^{(n)}}{\|\mathbf{v}_{s,[k]}^{(n)}\|}, \frac{\mathbf{v}_{t,[k]}^{(n)}}{\|\mathbf{v}_{t,[k]}^{(n)}\|} \right\rangle \right\rangle. \quad (23)$$

By Corollary 4.5, we can terminate the outer iteration when $\theta_{[k]} < \epsilon_{\text{outer}}$, which will be specified later. The maximum number of outer iterations is set to be 25.

5.2 Influence of Stopping Tolerances

We test different settings of tolerances: $\epsilon_{\text{inner}} = 10^{-3}, 10^{-4}, 10^{-5}$ and $\epsilon_{\text{outer}} = 10^{-3}, 10^{-4}, 10^{-5}$. The results are shown in Table 2, which are averaged over 10 trials for each case.

From Table 2, we can find that OD-ALM has a good performance on convergence: the outer iteration numbers are at most 12 on average for all cases. The running time would increase if we choose a smaller tolerance, but there is no improvement on the relative error for almost all cases. Therefore, we do not recommend using a too small tolerance in practical applications. We will use $\epsilon_{\text{inner}} = 10^{-4}, \epsilon_{\text{outer}} = 10^{-4}$ for synthetic tensors and $\epsilon_{\text{inner}} = 10^{-3}, \epsilon_{\text{outer}} = 10^{-3}$ for real-world tensors in all remaining tests.

5.3 Convergence Behaviour

We show the value of $\theta_{[k]}$ defined in (23), the relative change between successive outer iterates $\|\mathbf{v}_{[k]} - \mathbf{v}_{[k-1]}\|/\|\mathbf{v}_{[k-1]}\|$ and the number of inner iterations corresponding to each outer iteration in Figure 1 and Figure 2.

The value of $\theta_{[k]}$ is decreasing as k increases, but the situations differ greatly for different tensors. For example, $\theta_{[k]}$ of \mathcal{A}_7 is almost unchanged for the first five outer iterations, while $\theta_{[k]}$ of \mathcal{A}_6 decreases from more than 0.6 to less than 0.1 in the first five outer iterations. Usually, a big number of inner iterations brings a relatively big change of $\theta_{[k]}$. For example, for \mathcal{A}_3 , the number of inner iterations corresponding to $k = 2$ is more than 250, resulting in the difference between $\theta_{[1]}$ and $\theta_{[2]}$ being more than 0.4.

The relative change between successive outer iterates can be relatively big for some tensors even when k is big, e.g., \mathcal{A}_6 and \mathcal{A}_7 . For all cases, the relative change is relatively small between the last two outer iterates. The number of inner iterations reflects the relative change: A big number of inner iterations often results in a big relative change between successive outer iterates.

⁵ A Matlab implementation, adapted by Dianne P. O'Leary, is available at <http://www.cs.umd.edu/users/oleary/software/>.

Table 2 Results of OD-ALM under different stopping tolerances. Here “iter” is the number of outer iterations; the running time includes the time for Algorithms 1 and 2 and is measured in seconds

	ϵ_{outer}	ϵ_{inner}	\mathcal{A}_1	\mathcal{A}_2	\mathcal{A}_3	\mathcal{A}_4	\mathcal{A}_5	\mathcal{A}_6	\mathcal{A}_7	\mathcal{A}_8
Iter.	10^{-3}	10^{-3}	10	10	9	11	8	8	9	6
		10^{-4}	10	10	6	10	8	8	9	6
		10^{-5}	10	10	6	8	8	8	9	6
	10^{-4}	10^{-3}	11	11	10	12	9	9	9	7
		10^{-4}	11	11	7	11	8	9	9	7
		10^{-5}	11	11	6	9	8	9	9	7
	10^{-5}	10^{-3}	11	11	11	12	9	11	9	7
		10^{-4}	12	11	7	11	9	9	9	7
		10^{-5}	11	11	9	11	9	9	9	7
Time	10^{-3}	10^{-3}	1.1	1.0	1.3	0.6	4.8	13.2	15.8	15.3
		10^{-4}	2.6	1.3	3.2	0.5	13.9	24.4	19.6	21.8
		10^{-5}	4.7	1.8	4.6	0.4	24.5	34.6	22.9	30.0
	10^{-4}	10^{-3}	1.2	1.1	1.4	0.7	5.4	15.0	15.8	16.2
		10^{-4}	2.7	1.6	3.3	0.5	14.9	25.3	19.6	23.7
		10^{-5}	4.9	2.6	4.6	0.5	24.3	43.1	22.9	33.7
	10^{-5}	10^{-3}	1.2	1.1	1.6	0.7	4.9	15.7	15.6	16.2
		10^{-4}	2.7	1.6	2.9	0.6	15.3	26.1	19.6	23.8
		10^{-5}	4.9	3.9	5.8	0.9	24.5	41.9	23.3	33.6
RErr	10^{-3}	10^{-3}	0.9954	0.0559	0.0640	0.0994	0.1831	0.2379	0.2931	0.2278
		10^{-4}	0.9954	0.0559	0.0267	0.0994	0.1831	0.2378	0.2931	0.2278
		10^{-5}	0.9954	0.0559	0.0245	0.0993	0.1831	0.2378	0.2931	0.2278
	10^{-4}	10^{-3}	0.9954	0.0559	0.0640	0.0994	0.1831	0.2379	0.2931	0.2278
		10^{-4}	0.9954	0.0559	0.0227	0.0994	0.1831	0.2378	0.2931	0.2278
		10^{-5}	0.9954	0.0559	0.0245	0.0993	0.1831	0.2378	0.2931	0.2278
	10^{-5}	10^{-3}	0.9954	0.0559	0.0640	0.0994	0.1831	0.2379	0.2931	0.2278
		10^{-4}	0.9954	0.0559	0.0227	0.0994	0.1831	0.2378	0.2931	0.2278
		10^{-5}	0.9954	0.0559	0.0245	0.0993	0.1831	0.2378	0.2931	0.2278

5.4 Comparison with Other Methods

We compare our method with CP-ALS, the low rank orthogonal approximation of tensors (LROAT) [5] and the high-order power method for orthogonal low rank decomposition (OLRD-HOP) [34]. The method CP-ALS fits a CP decomposition (6). The method LROAT fits an $(1, \dots, N)$ -orthogonal decomposition, and OLRD-HOP fits an (N) -orthogonal decomposition. CP-ALS, LROAT and OLRD-HOP are all with the truncated HOSVD initialization. CP-ALS terminates if the relative change in the function value is less than 10^{-8} . LROAT and OLRD-HOP terminate if the relative change between successive iterates is less than 10^{-8} . The maximum number of iterations is set to be 500 for all these three methods.

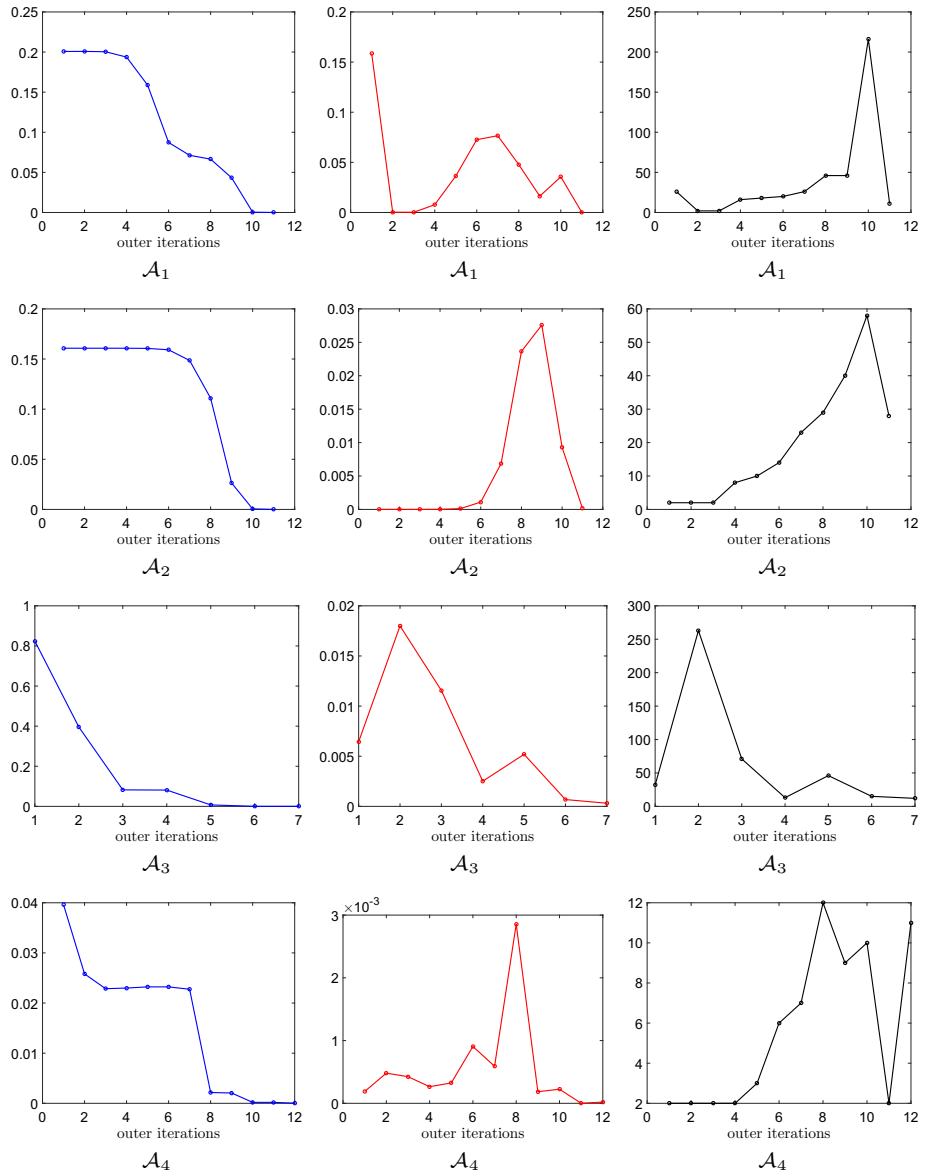


Fig. 1 The convergence behaviour of OD-ALM on $\mathcal{A}_1, \dots, \mathcal{A}_4$. The first column is about $\theta_{[k]}$, the second column is about $\frac{\|v_{[k]} - v_{[k-1]}\|}{\|v_{[k-1]}\|}$, and the last column is about the number of inner iterations. All values are shown as functions of the number of outer iterations

The results of the running time and the relative error are shown in Table 3, which are averaged over 10 trials for each case.

We can see that our method is much slower than the other methods. As discussed in [1], the time cost of one outer iteration of OD-ALM is of the same order of magnitude with

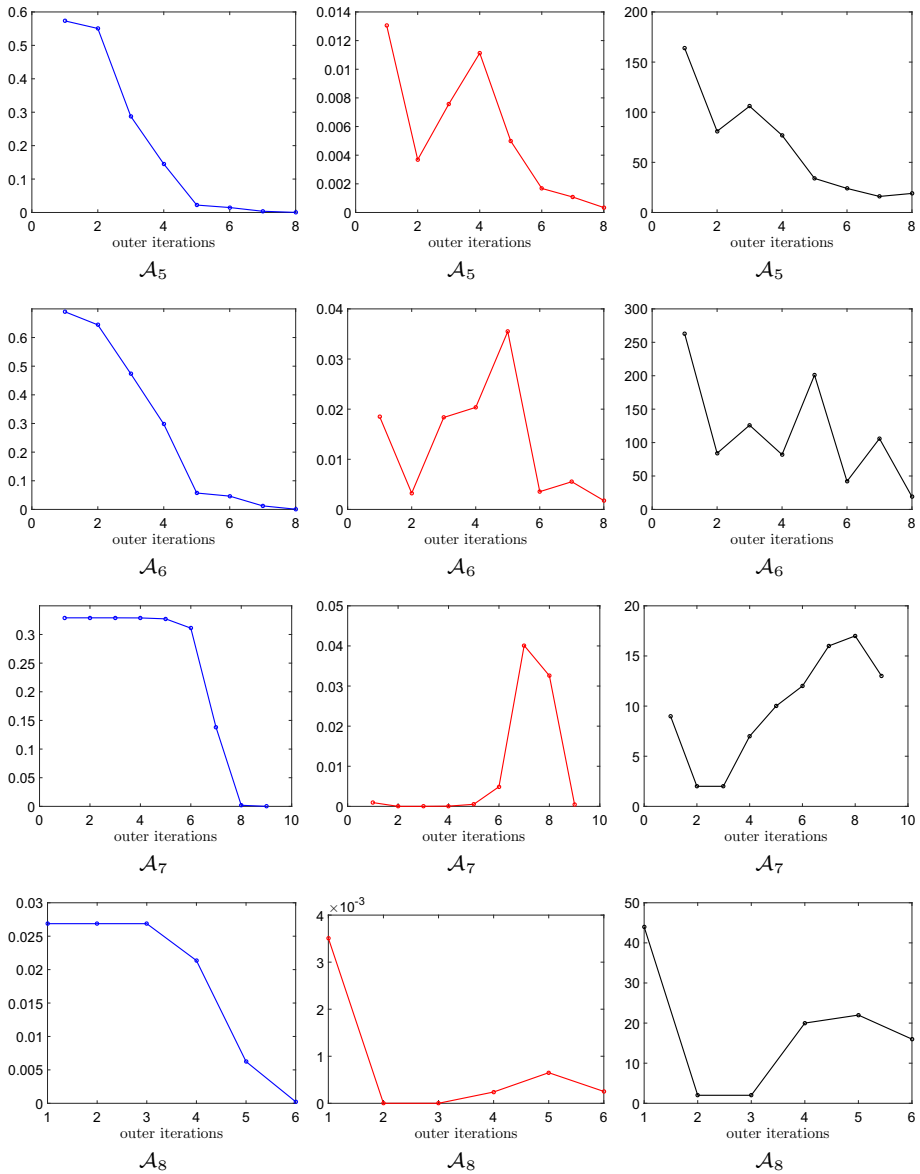


Fig. 2 The convergence behaviour of OD-ALM on $\mathcal{A}_5, \dots, \mathcal{A}_8$. The three columns have the same meaning as in Figure 1

CP-ALS. OD-ALM needs several outer iterations, resulting in a much longer time cost than CP-ALS. The running time of LROAT and OLRD-HOP is close to that of CP-ALS.

As for the relative error, CP-ALS is the best, OD-ALM is the second best, and OLRD-HOP outperforms LROAT. This is not surprising because of the relationships among the decompositions fitted by different methods. For \mathcal{A}_4 whose ground truth is an orthogonal rank-5 tensor, the OD-ALM RErr is less than the noise level 0.1, which demonstrates the

Table 3 Comparison results of different methods, where OD-ALM has been combined with Algorithm 2. The running time is measured in seconds

	Method	\mathcal{A}_1	\mathcal{A}_2	\mathcal{A}_3	\mathcal{A}_4	\mathcal{A}_5	\mathcal{A}_6	\mathcal{A}_7	\mathcal{A}_8
Time	CP-ALS	0.3	0.1	0.8	0.1	1.3	1.6	1.7	5.1
	OD-ALM	2.7	1.6	3.3	0.5	4.8	13.2	15.8	15.3
	LROAT	2.2	0.07	0.06	0.06	0.7	1.3	3.8	8.4
	OLRD-HOP	0.6	0.07	1.3	1.3	2.1	2.5	1.2	2.9
RErr	CP-ALS	0.9953	0	0.0070	0.0993	0.1822	0.2363	0.2857	0.2278
	OD-ALM	0.9954	0.0559	0.0227	0.0994	0.1831	0.2379	0.2931	0.2278
	LROAT	0.9957	0.2890	0.1728	0.1640	0.3504	0.3263	0.4513	0.2530
	OLRD-HOP	0.9954	0.1604	0.1117	0.1478	0.3333	0.3174	0.4510	0.2525

effectiveness of our method. In addition, we can find that the difference between the CP-ALS RErr and the OD-ALM RErr is very small for real-world tensors. For \mathcal{A}_8 , the results of these two methods are even the same. This suggests the potential of orthogonal decompositions in fitting real-world tensors. The small gap between the CP-ALS RErr and the OD-ALM RErr also indicates the effectiveness of our method in some sense.

Suppose $\mathbf{U}_j^{(n)}$ is the n th normalized factor matrix corresponding to the final result for \mathcal{A}_j obtained by our method. We record the results of $\mathbf{U}_j^{(n)\top} \mathbf{U}_j^{(n)}$ for $j = 3, 5$ in one running:

$$\mathbf{U}_3^{(1)\top} \mathbf{U}_3^{(1)} = \begin{bmatrix} 1 & 0.6089 & 0.6264 & -0.3196 & 0 \\ 0.6089 & 1 & 0.9814 & 0.5454 & 0.7771 \\ 0.6264 & 0.9814 & 1 & 0.4745 & 0.7039 \\ -0.3196 & 0.5454 & 0.4745 & 1 & 0.9472 \\ 0 & 0.7771 & 0.7039 & 0.9472 & 1 \end{bmatrix}$$

$$\mathbf{U}_3^{(2)\top} \mathbf{U}_3^{(2)} = \begin{bmatrix} 1 & 0 & -0.1713 & -0.9277 & -0.8513 \\ 0 & 1 & 0.9685 & 0.3720 & 0.5199 \\ -0.1713 & 0.9685 & 1 & 0.5136 & 0.6367 \\ -0.9277 & 0.3720 & 0.5136 & 1 & 0.9853 \\ -0.8513 & 0.5199 & 0.6367 & 0.9853 & 1 \end{bmatrix}$$

$$\mathbf{U}_3^{(3)\top} \mathbf{U}_3^{(3)} = \begin{bmatrix} 1 & 0.2055 & -0.5054 & -0.9921 & -0.9775 \\ 0.2055 & 1 & 0.7289 & -0.0832 & 0 \\ -0.5054 & 0.7289 & 1 & 0.6030 & 0.6618 \\ -0.9921 & -0.0832 & 0.6030 & 1 & 0.9962 \\ -0.9775 & 0 & 0.6618 & 0.9962 & 1 \end{bmatrix}$$

$$\mathbf{U}_3^{(4)\top} \mathbf{U}_3^{(4)} = \begin{bmatrix} 1 & -1 & 0 & 0 & -0.9996 \\ -1 & 1 & 0 & 0 & 0.9996 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ -0.9996 & 0.9996 & 0 & 0 & 1 \end{bmatrix};$$

$$\mathbf{U}_5^{(1)\top} \mathbf{U}_5^{(1)} = \begin{bmatrix} 1 & 0 & 0.7831 & -0.4958 & 0 \\ 0 & 1 & 0.0954 & 0.0805 & -0.3413 \\ 0.7831 & 0.0954 & 1 & 0 & 0 \\ -0.4958 & 0.0805 & 0 & 1 & -0.2793 \\ 0 & -0.3413 & 0 & -0.2793 & 1 \end{bmatrix}$$

$$\mathbf{U}_5^{(2)\top} \mathbf{U}_5^{(2)} = \begin{bmatrix} 1 & 0.7868 & 0 & 0 & -0.1452 \\ 0.7868 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & -0.6186 & -0.0751 \\ 0 & 0 & -0.6186 & 1 & 0 \\ -0.1452 & 0 & -0.0751 & 0 & 1 \end{bmatrix}$$

$$\mathbf{U}_5^{(3)\top} \mathbf{U}_5^{(3)} = \begin{bmatrix} 1 & 0.9091 & 0.9243 & 0.9867 & -0.9640 \\ 0.9091 & 1 & 0.9992 & 0.9629 & -0.9864 \\ 0.9243 & 0.9992 & 1 & 0.9720 & -0.9920 \\ 0.9867 & 0.9629 & 0.9720 & 1 & -0.9933 \\ -0.9640 & -0.9864 & -0.9920 & -0.9933 & 1 \end{bmatrix}.$$

We also compute $\mathbf{U}_j^{(n)\top} \mathbf{U}_j^{(n)}$ for other tensors and find that the appearance of zeros in $\mathbf{U}_j^{(n)\top} \mathbf{U}_j^{(n)}$ is not regular. Therefore, strongly orthogonal decompositions cannot replace orthogonal decompositions in practical

6 Conclusion

We present several properties of orthogonal rank. Orthogonal rank is different from tensor rank in many aspects. For example, a subtensor may have a larger orthogonal rank than the whole tensor, and orthogonal rank is lower semicontinuous.

To tackle the complicated orthogonality constraints, we employ the augmented Lagrangian method to convert the original problem into an unconstrained problem. The gradient of the objective function has a good structure, inspiring us to use gradient-based optimization methods to solve each subproblem. A novel orthogonalization process is developed to make the final result satisfy the orthogonality condition exactly. Numerical experiments show that the proposed method has a great advantage over the existing methods for strongly orthogonal decompositions in terms of the approximation error.

The main drawback of our method is the time cost. This is because the time cost of one outer iteration of OD-ALM is of the same order of magnitude with that of CP-ALS, which is not very short, and we need several outer iterations to obtain the final result. Although the ill-conditioning is not so severe for the augmented Lagrangian method compared to the penalty method, preconditioning is a possible way to speed up. For preconditioning of optimization methods for CP decompositions, one can refer to [10, 32]. Preconditioning for OD-ALM can be studied as future work. A better strategy is to design an algorithm with a framework different from the augmented Lagrangian method. This may need further exploration of orthogonal decompositions.

Acknowledgements The author is extremely grateful to the two anonymous referees for their valuable feedback, which improved this paper significantly. This work was partially supported by the National Natural Science Foundation of China (12201319).

Author Contributions CZ is the single author of the manuscript and responsible for this work.

Funding This work was partially supported by the National Natural Science Foundation of China (12201319).

Data Availability The datasets analysed during the current study are public available and we have provided the URLs when using them.

Declarations

Conflict of interest The author declares he/she has no financial interest.

References

1. Acar, E., Dunlavy, D.M., Kolda, T.G.: A scalable optimization approach for fitting canonical tensor decompositions. *J. Chemom.* **25**(2), 67–86 (2011)
2. Bader, B.W., Kolda, T.G. et al.: MATLAB Tensor Toolbox Version 3.0-dev. Available online, Oct. (2017)
3. Bertsekas, D.P.: *Constrained Optimization and Lagrange Multiplier Methods*. Academic press, Cambridge (1982)
4. Carroll, J.D., Chang, J.-J.: Analysis of individual differences in multidimensional scaling via an N-way generalization of “Eckart-Young” decomposition. *Psychometrika* **35**(3), 283–319 (1970)

5. Chen, J., Saad, Y.: On the tensor SVD and the optimal low rank orthogonal approximation of tensors. *SIAM J. Matrix Anal. Appl.* **30**(4), 1709–1734 (2008)
6. Comon, P.: Independent component analysis, A new concept? *Signal Process.* **36**(3), 287–314 (1994)
7. Conn, A.R., Gould, N., Sartenaer, A., Toint, P.L.: Convergence properties of an augmented Lagrangian algorithm for optimization with a combination of general equality and linear constraints. *SIAM J. Optim.* **6**(3), 674–703 (1996)
8. De Lathauwer, L., De Moor, B., Vandewalle, J.: A multilinear singular value decomposition. *SIAM J. Matrix Anal. Appl.* **21**(4), 1253–1278 (2000)
9. De Silva, V., Lim, L.-H.: Tensor rank and the ill-posedness of the best low-rank approximation problem. *SIAM J. Matrix Anal. Appl.* **30**(3), 1084–1127 (2008)
10. De Sterck, H., Howse, A.J.: Nonlinearly preconditioned L-BFGS as an acceleration mechanism for alternating least squares with application to tensor decomposition. *Num. Linear Algebra Appl.* **25**(6), e2202 (2018)
11. Eckart, C., Young, G.: The approximation of one matrix by another of lower rank. *Psychometrika* **1**(3), 211–218 (1936)
12. Espig, M., Hackbusch, W.: A regularized Newton method for the efficient approximation of tensors represented in the canonical tensor format. *Numer. Math.* **122**(3), 489–525 (2012)
13. Guan, Y., Chu, D.: Numerical computation for orthogonal low-rank approximation of tensors. *SIAM J. Matrix Anal. Appl.* **40**(3), 1047–1065 (2019)
14. Harshman, R.A. et al.: Foundations of the PARAFAC procedure: models and conditions for an “explanatory” multimodal factor analysis. (1970)
15. Hästad, J.: Tensor rank is NP-complete. *J. Algorithms* **11**(4), 644–654 (1990)
16. Hillar, C.J., Lim, L.-H.: Most tensor problems are NP-hard. *J. ACM (JACM)* **60**(6), 45 (2013)
17. Kolda, T.G.: Orthogonal tensor decompositions. *SIAM J. Matrix Anal. Appl.* **23**(1), 243–255 (2001)
18. Kolda, T.G.: A counterexample to the possibility of an extension of the Eckart-Young low-rank approximation theorem for the orthogonal rank tensor decomposition. *SIAM J. Matrix Anal. Appl.* **24**(3), 762–767 (2003)
19. Kolda, T.G., Bader, B.W.: Tensor decompositions and applications. *SIAM Rev.* **51**(3), 455–500 (2009)
20. Krijnen, W.P., Dijkstra, T.K., Stegeman, A.: On the non-existence of optimal solutions and the occurrence of “degeneracy” in the CANDECOMP/PARAFAC model. *Psychometrika* **73**(3), 431–439 (2008)
21. Kruskal, J.B.: Three-way arrays: rank and uniqueness of trilinear decompositions, with application to arithmetic complexity and statistics. *Linear Algebra Appl.* **18**(2), 95–138 (1977)
22. Li, Z., Nakatsukasa, Y., Soma, T., Uschmajew, A.: On orthogonal tensors and best rank-one approximation ratio. *SIAM J. Matrix Anal. Appl.* **39**(1), 400–425 (2018)
23. Lim, L.-H., Comon, P.: Blind multilinear identification. *IEEE Trans. Inf. Theory* **60**(2), 1260–1280 (2013)
24. Martin, C.D.M., Van Loan, C.F.: A Jacobi-type method for computing orthogonal tensor decompositions. *SIAM J. Matrix Anal. Appl.* **30**(3), 1219–1232 (2008)
25. More, J.J., Thuente, D.J.: Line search algorithms with guaranteed sufficient decrease. *ACM Trans. Math. Softw.* **20**(3), 286–307 (1994)
26. Nazih, M., Minaoui, K., Comon, P.: Using the proximal gradient and the accelerated proximal gradient as a canonical polyadic tensor decomposition algorithms in difficult situations. *Signal Process.* **171**, 107472 (2020)
27. Nocedal, J., Wright, S.: Numerical Optimization. Springer Science & Business Media, New York (2006)
28. Rockafellar, R.T., Wets, R.J.-B.: Variational Analysis. Springer Science & Business Media, New York (2009)
29. Seeling, P., Reisslein, M.: Video transport evaluation with H. 264 video traces. *IEEE Commun. Surv. Tutor.* **14**(4), 1142–1165 (2011)
30. Sidiropoulos, N.D., Bro, R.: On the uniqueness of multilinear decomposition of N-way arrays. *J. Chemometr. J. Chemometr. Soc.* **14**(3), 229–239 (2000)
31. Sørensen, M., De Lathauwer, L., Comon, P., Icart, S., Deneire, L.: Canonical polyadic decomposition with a columnwise orthonormal factor matrix. *SIAM J. Matrix Anal. Appl.* **33**(4), 1190–1213 (2012)
32. Sterck, H.D.: A nonlinear GMRES optimization algorithm for canonical tensor decomposition. *SIAM J. Sci. Comput.* **34**(3), A1351–A1379 (2012)
33. Sun, W., Yuan, Y.-X.: Optimization Theory and Methods: Nonlinear Programming. Springer Optimization and Its Applications. Springer Science & Business Media, New York (2010)
34. Wang, L., Chu, M.T., Yu, B.: Orthogonal low rank tensor approximation: alternating least squares method and its global convergence. *SIAM J. Matrix Anal. Appl.* **36**(1), 1–19 (2015)
35. Yang, Y.: The epsilon-alternating least squares for orthogonal low-rank tensor approximation and its global convergence. *SIAM J. Matrix Anal. Appl.* **41**(4), 1797–1825 (2020)

36. Zhu, F., Wang, Y., Fan, B., Xiang, S., Meng, G., Pan, C.: Spectral unmixing via data-guided sparsity. *IEEE Trans. Image Process.* **23**(12), 5412–5427 (2014)

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.