# ON THE TENSOR SVD AND THE OPTIMAL LOW RANK ORTHOGONAL APPROXIMATION OF TENSORS[*]

JIE CHEN[†] AND YOUSEF SAAD[†]

**Abstract.** It is known that a higher order tensor does not necessarily have an optimal low rank approximation, and that a tensor might not be orthogonally decomposable (i.e., admit a tensor SVD). We provide several sufficient conditions which lead to the failure of the tensor SVD, and characterize the existence of the tensor SVD with respect to the Higher Order SVD (HOSVD). In face of these difficulties to generalize standard results known in the matrix case to tensors, we consider the low rank orthogonal approximation of tensors. The existence of an optimal approximation is theoretically guaranteed under certain conditions, and this optimal approximation yields a tensor decomposition where the diagonal of the core is maximized. We present an algorithm to compute this approximation and analyze its convergence behavior. Numerical experiments indicate a linear convergence rate for this algorithm.

**Key words.** multilinear algebra, singular value decomposition, tensor decomposition, low rank approximation

**AMS subject classifications.** 15A69, 15A18

**1. Introduction.** There has been renewed interest in studying the properties and decompositions of tensors (also known as $N$-way arrays or multidimensional arrays) in numerical linear algebra in recent years [30, 13, 12, 43, 17, 9, 28, 29, 15, 11]. The tensor approximation techniques have been fruitfully applied in various areas which include among others, chemometrics [38, 4], signal processing [10, 8], vision and graphics [41, 42], and network analysis [31, 1]. From the point of view of practical applications, the matrix SVD and the optimal rank-$r$ approximation of matrices (a.k.a. Eckart-Young theorem [18]) are of particular interest, and it would be nice if these properties could be directly generalized to higher order tensors. However, for any order $N \geq 3$, de Silva and Lim [17] showed that the problem of optimal low rank approximation of higher order tensors is ill-posed for many ranks $r$, and that this ill-posedness is not rare for order-3 tensors. Furthermore, Kolda presented numerous examples to illustrate the difficulties of orthogonal tensor decompositions [28, 29]. These studies revealed many aspects of the dissimilarities between tensors and matrices, in spite of the fact that higher order tensors are multidimensional generalizations of matrices.

The most commonly used generalization of the matrix SVD to higher order tensors to date is the *Higher Order Singular Value Decomposition* (HOSVD) [12]. The HOSVD decomposes an order-$N$ tensor into a core tensor that has the same size as the original tensor together with $N$ orthogonal[1] side-matrices. Although this decomposition preserves many nice aspects of the matrix SVD (e.g., the core has the all-orthogonality property and the ordering property), a notable difference is that the core is in general not diagonal. Hence, in contrast with the matrix SVD, the HOSVD cannot be written as a sum of a few orthogonal outer-product terms[2].

---

[†]Department of Computer Science and Engineering, University of Minnesota at Twin Cities, MN 55455. Email: {jchen, saad}@cs.umn.edu.

[1]Throughout this paper, a matrix $A \in \mathbb{R}^{m \times n}$, $m \geq n$, is said to be *orthogonal* if $A^T A = I$. This generalizes the definition for square matrices.

[2]For discussions of orthogonality, see Section 2.4.

There exist three well-known approximations to higher order tensors: (1) the rank-1 approximation [13, 43, 27]; (2) the rank-$(r_1, r_2, \ldots, r_N)$ approximation with a full core and $N$ orthogonal side-matrices (in the *Tucker/HOOI* fashion) [40, 13]; and (3) the approximation using $r$ outer-product terms (in the *CANDECOMP/PARAFAC* fashion) [6, 19]. Note that the approximated tensor in case (3) might have rank less than $r$. Among these approximations, the rank-1 approximation [17] and the rank-$(r_1, r_2, \ldots, r_N)$ approximation are theoretically guaranteed to have a global optimum. In practical applications, the three approximations are generally computed using an alternating least squares (ALS) method [33, 3, 25] (the so-called "workhorse" algorithm [30]), although many other methods have also been proposed [34, 37, 43, 28, 15, 11]. The convergence behavior of the ALS method is theoretically unknown except under a few strong conditions [32]. Besides, it has long been observed that the ALS method for the PARAFAC model may converge extremely slowly if at all [36, 26]. An illustration of this phenomenon is given in the Appendix.

Kolda [28] investigated several orthogonal decompositions of tensors related to different definitions of orthogonality, including *orthogonal rank decomposition*, *complete orthogonal rank decomposition* and *strong orthogonal rank decomposition*. These decompositions might not be unique, or even exist. Among these definitions, only the *complete orthogonality* gives a situation which parallels that of the matrix SVD. This approach demands that the side-matrices all be orthogonal, in which case we use the term *tensor singular value decomposition* (tensor SVD, see Definition 4.1) in this paper. Zhang and Golub [43] proved that for all tensors of order $N \geq 3$, the tensor SVD is unique (up to signs) if it exists, and that the incremental rank-1 approximation approach will compute this decomposition.

The following contributions are made in this paper:

1. Sufficient conditions indicating which tensors fail to have a tensor SVD are given. These conditions are related to the rank, the order, and the dimensions of the tensor, and hence can be viewed as generalizations of results given in the literature with specific examples. Furthermore, the existence of the tensor SVD can be characterized by the diagonality of the core in the HOSVD of the tensor.

2. A form of low rank approximations—one that requires a diagonal core and orthogonal side-matrices—is discussed. Theoretically the global optimum of this approximation can be attained for any (appropriate) rank. We present an iterative algorithm to compute this approximation and analyze its convergence behavior.

3. The proposed approximation at the maximally possible rank leads to a decomposition of the tensor, where the diagonal of the core is maximized. This "maximal diagonality" for symmetric order-3 [14] and 4 [7] tensors and for general order-3 tensors [15, 24, 35] has been previously investigated and Jacobi algorithms were used in the cited papers, but our discussion is in a more general context and the proposed algorithm is not of a Jacobi type.

**2. Tensor algebra.** In this section, we briefly review some concepts and notions that are used throughout the paper. A *tensor* is a multidimensional array of data whose elements are referred by using multiple indices. The number of indices required is called the *order* of a tensor. We use

$$\mathcal{A} = (a_{i_1, i_2, \ldots, i_N}) \in \mathbb{R}^{d_1 \times d_2 \times \cdots \times d_N}$$

to denote a tensor $\mathcal{A}$ of order $N$. For $n = 1, 2, \ldots, N$, $d_n$ is the $n$-th *dimension* of $\mathcal{A}$. As a special case, a vector is an order-1 tensor and a matrix is an order-2 tensor.

**2.1. Unfoldings and mode-$n$ products.** It is hard to visualize tensors of order $N > 3$. They can be flexibly represented when "unfolded" into matrices. The *unfolding* of a tensor along *mode $n$* is a matrix of dimension $d_n \times (d_{n+1} \cdots d_N d_1 \cdots d_{n-1})$. We denote the mode-$n$ unfolding of tensor $\mathcal{A}$ by $A_{(n)}$. Each column of $A_{(n)}$ is a column of $\mathcal{A}$ along the $n$-th mode.

An important operation for a tensor is the *tensor-matrix multiplication*, also known as *mode-$n$ product*. Given a tensor $\mathcal{A} \in \mathbb{R}^{d_1 \times d_2 \times \cdots \times d_N}$ and a matrix $M \in \mathbb{R}^{c_n \times d_n}$, the mode-$n$ product is a tensor

$$\mathcal{B} = \mathcal{A} \times_n M \in \mathbb{R}^{d_1 \times \cdots \times d_{n-1} \times c_n \times d_{n+1} \cdots \times d_N}$$

where

$$b_{i_1,\ldots,i_{n-1},j_n,i_{n+1},\ldots,i_N} := \sum_{i_n=1}^{d_n} a_{i_1,\ldots,i_{n-1},i_n,i_{n+1},\ldots,i_N} m_{j_n,i_n}$$

for $j_n = 1, 2, \ldots, c_n$. In matrix representation, this is

$$B_{(n)} = M A_{(n)}. \tag{2.1}$$

**2.2. Inner products and tensor norms.** The *inner product* of two tensors $\mathcal{A}$ and $\mathcal{B}$ of the same size is defined by

$$\langle \mathcal{A}, \mathcal{B} \rangle_F := \sum_{i_N=1}^{d_N} \cdots \sum_{i_1=1}^{d_1} a_{i_1,\ldots,i_N} b_{i_1,\ldots,i_N}.$$

and the *norm* induced from this inner product is

$$\|\mathcal{A}\|_F := \sqrt{\langle \mathcal{A}, \mathcal{A} \rangle_F}.$$

We say that $\mathcal{A}$ is a *unit tensor* if $\|\mathcal{A}\|_F = 1$. When $N = 2$, $\|\mathcal{A}\|_F$ is the Frobenius norm of matrix $\mathcal{A}$. The norm of a tensor is equal to the Frobenius norm of the unfolding of the tensor along any mode:

$$\|\mathcal{A}\|_F = \|A_{(n)}\|_F, \quad \text{for } n = 1, \ldots, N.$$

**2.3. Tensor products and outer products of vectors.** The *tensor product* of an order-$N$ tensor $\mathcal{A} \in \mathbb{R}^{d_1 \times d_2 \times \cdots \times d_N}$ and an order-$N'$ tensor $\mathcal{B} \in \mathbb{R}^{c_1 \times c_2 \times \cdots \times c_{N'}}$ is an order-$(N + N')$ tensor

$$\mathcal{C} = \mathcal{A} \otimes \mathcal{B} \in \mathbb{R}^{d_1 \times \cdots \times d_N \times c_1 \times \cdots \times c_{N'}},$$

where

$$c_{i_1,\ldots,i_N,j_1,\ldots,j_{N'}} := a_{i_1,\ldots,i_N} b_{j_1,\ldots,j_{N'}}.$$

Note that the operator $\otimes$ for tensor products unfortunately coincides with the one used to denote the Kronecker product of two matrices. In particular, the tensor product of two matrices (order-2 tensors) is an order-4 tensor, while the Kronecker product of two matrices is again a matrix. The reader shall not be confused by this notation since in this paper Kronecker products are not involved.

The *outer product* of $N$ (column) vectors, which generalizes the standard outer product of two vectors (a rank-1 matrix), is a special case of tensor products. The outer product of $N$ vectors $x^{(n)} \in \mathbb{R}^{d_n}$, is an order-$N$ tensor of dimension $d_1 \times d_2 \times \cdots \times d_N$:

$$\mathcal{X} = x^{(1)} \otimes x^{(2)} \otimes \cdots \otimes x^{(N)}.$$

The $(i_1, i_2, \ldots, i_N)$-entry of $\mathcal{X}$ is $\prod_{n=1}^N (x^{(n)})_{i_n}$, where $(x^{(n)})_{i_n}$ denotes the $i_n$-th entry of vector $x^{(n)}$. The tensor $\mathcal{X}$ is also called a *rank-1 tensor*. The *rank* of a tensor is defined in Section 3.

It can be verified that the mode-$n$ product of a rank-1 tensor $\mathcal{X}$ with a matrix $M$ can be computed as follows:

$$\mathcal{X} \times_n M = x^{(1)} \otimes \cdots \otimes x^{(n-1)} \otimes \left( M x^{(n)} \right) \otimes x^{(n+1)} \cdots \otimes x^{(N)},$$

and that the inner product of $\mathcal{X}$ with a general tensor $\mathcal{A}$ is

$$
\begin{aligned}
\langle \mathcal{A}, \mathcal{X} \rangle_F &= \left\langle \mathcal{A}, x^{(1)} \otimes x^{(2)} \otimes \cdots \otimes x^{(N)} \right\rangle_F \\
&= \mathcal{A} \times_1 {x^{(1)}}^T \times_2 {x^{(2)}}^T \times \cdots \times_N {x^{(N)}}^T.
\end{aligned}
$$

If $\mathcal{U} = u^{(1)} \otimes u^{(2)} \otimes \cdots \otimes u^{(N)}$ and $\mathcal{V} = v^{(1)} \otimes v^{(2)} \otimes \cdots \otimes v^{(N)}$ are two rank-1 tensors then

$$\langle \mathcal{U}, \mathcal{V} \rangle_F = \prod_{n=1}^N \left\langle u^{(n)}, v^{(n)} \right\rangle,$$

where $\langle \cdot, \cdot \rangle$ denotes the standard Euclidean inner product of two vectors. A consequence of the above relation is that $\|\mathcal{U}\|_F$ is the product of the 2-norms of the vectors $u^{(n)}$'s.

**2.4. Orthogonality of tensors.** Two tensors $\mathcal{A}$ and $\mathcal{B}$ of the same size are *F-orthogonal* (*Frobenius orthogonal*) if their inner product is zero, i.e.,

$$\langle \mathcal{A}, \mathcal{B} \rangle_F = 0.$$

For rank-1 tensors $\mathcal{U} = u^{(1)} \otimes u^{(2)} \otimes \cdots \otimes u^{(N)}$ and $\mathcal{V} = v^{(1)} \otimes v^{(2)} \otimes \cdots \otimes v^{(N)}$, the above definition implies that they are F-orthogonal if

$$\prod_{n=1}^N \left\langle u^{(n)}, v^{(n)} \right\rangle = 0.$$

This leads to other options for defining orthogonality for two rank-1 tensors. The paper [28] discussed two cases:

1. Complete orthogonality: $\left\langle u^{(n)}, v^{(n)} \right\rangle = 0$ for all $n = 1, \ldots, N$.

2. Strong orthogonality: For all $n$, either $\left\langle u^{(n)}, v^{(n)} \right\rangle = 0$ or $u^{(n)}$ and $v^{(n)}$ are collinear, but there is at least one $\ell$ such that $\left\langle u^{(\ell)}, v^{(\ell)} \right\rangle = 0$.

In this paper we will simply use the term *orthogonal* for two outer products that are completely orthogonal (case 1).
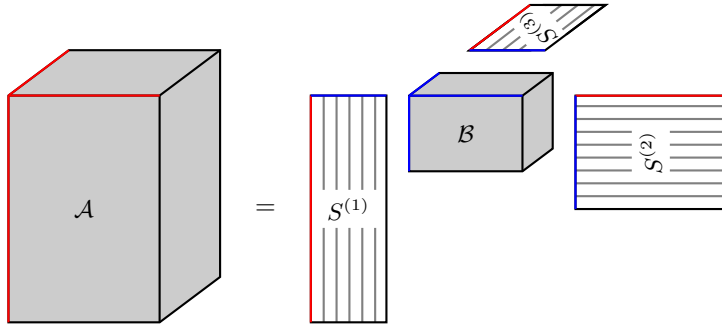
FIG. 2.1. *A decomposition of an order-3 tensor $\mathcal{A}$ as $\mathcal{B} \times_1 S^{(1)} \times_2 S^{(2)} \times_3 S^{(3)}$.*

**2.5. Tensor decompositions.** A *decomposition* of a tensor $\mathcal{A} \in \mathbb{R}^{d_1 \times d_2 \times \cdots \times d_N}$ is of the form

$$\mathcal{A} = \mathcal{B} \times_1 S^{(1)} \times_2 S^{(2)} \times \cdots \times_N S^{(N)},$$

where $\mathcal{B} \in \mathbb{R}^{c_1 \times c_2 \times \cdots \times c_N}$ is called the *core tensor*, and $S^{(n)} \in \mathbb{R}^{d_n \times c_n}$ for $n = 1, \ldots, N$ are called *side-matrices*. An illustration is given in Figure 2.1.

Let $s_i^{(n)}$ be the $i$-th column of $S^{(n)}$. The decomposition of $\mathcal{A}$ can equivalently be written as a linear combination of rank-1 tensors:

$$\mathcal{A} = \sum_{i_N=1}^{c_N} \cdots \sum_{i_1=1}^{c_1} b_{i_1, i_2, \ldots, i_N} s_{i_1}^{(1)} \otimes s_{i_2}^{(2)} \otimes \cdots \otimes s_{i_N}^{(N)}. \tag{2.2}$$

In particular, if $\mathcal{B}$ is diagonal, i.e., $b_{i_1, i_2, \ldots, i_N} = 0$ except when $i_1 = i_2 = \cdots = i_N$, then

$$\mathcal{A} = \sum_{i=1}^{r} b_{ii\ldots i} s_i^{(1)} \otimes s_i^{(2)} \otimes \cdots \otimes s_i^{(N)} \tag{2.3}$$

where $r = \min\{c_1, \ldots, c_N\}$.

In the literature, the term "decomposition" is often used when "approximation" is meant instead. The *Tucker3 decomposition* is an approximation in the form of the right-hand side of (2.2), for given dimensions $c_1, c_2, \ldots, c_N$. Usually, it is required that $c_n$ is less than the rank of $A_{(n)}$ for all $n$, otherwise the problem is trivial. The HOOI approach computes this approximation with an additional property that all the $S^{(n)}$'s are orthogonal matrices. The *CANDECOMP/PARAFAC decomposition* is an approximation in the form of the right-hand side of (2.3), for a pre-specified $r$. Usually, $r$ is smaller than the smallest dimension of all modes of $\mathcal{A}$, although requiring a larger $r$ is also possible in the ALS and other algorithms. As will be discussed in the next section, the smallest $r$ that satisfies equality (2.3) is the rank of the tensor $\mathcal{A}$.

**3. Tensor ranks.** The rank of a tensor causes difficulties when attempting to generalize matrix properties to tensors. There are several possible generalizations of the notion of rank. The *n-rank* of a tensor $\mathcal{A} \in \mathbb{R}^{d_1 \times d_2 \times \cdots \times d_N}$, for $n = 1, \ldots, N$, denoted by $\mathrm{rank}_n(\mathcal{A})$, is the rank of the unfolding $A_{(n)}$:

$$\mathrm{rank}_n(\mathcal{A}) := \mathrm{rank}(A_{(n)}).$$

The *(outer-product) rank* of $\mathcal{A}$, denoted by $\text{rank}(\mathcal{A})$, is defined as

$$\text{rank}(\mathcal{A}) := \min\left\{ r \,\middle|\, \exists\, x_i^{(n)} \in \mathbb{R}^{d_n}, \;\; i = 1, \ldots, r, \;\; n = 1, \ldots, N, \right.$$

$$\left. \text{s.t. } \mathcal{A} = \sum_{i=1}^{r} x_i^{(1)} \otimes x_i^{(2)} \otimes \cdots \otimes x_i^{(N)} \right\}.$$

Hence, a tensor is the outer product of $N$ vectors if and only if it has rank one, and the rank of a general tensor $\mathcal{A}$ is the minimum number of rank-1 tensors that sum to $\mathcal{A}$.

There are a few notable differences between the notion of rank for matrices and that for tensors:

1. For $N = 2$, i.e., when $\mathcal{A}$ is a matrix, $\text{rank}_1(\mathcal{A})$ is the row rank, $\text{rank}_2(\mathcal{A})$ is the column rank, and $\text{rank}(\mathcal{A})$ is the outer-product rank, and they are all equal. However, for higher order tensors $(N > 2)$, in general, the $n$-ranks are different for different modes $n$, and they are different from $\text{rank}(\mathcal{A})$ [12]. Furthermore, the rank of a matrix $A$ can not be larger than the smallest dimension of both modes of $A$, but for tensors this is no longer true, i.e., the rank can be larger than the smallest dimension of the tensor [12].

2. The matrix SVD yields one possible way of writing a matrix as a sum of outer-product terms, and the number of nonzero singular values is equal to the rank of the matrix. However, a tensor SVD does not always exist (see Section 4), but if it indeed does, it is unique up to signs [34, 43] and the number of singular values is equal to the rank of the tensor (see Definition 4.1 and Proposition 4.2).

3. It is well-known that the optimal rank-$r$ approximation of a matrix is simply its truncated SVD. However some tensors may fail to have an optimal rank-$r$ approximation [17]. If such an approximation exists, it is unclear whether it can be written in the form of a diagonal core multiplied by orthogonal side-matrices.

Next are some lemmas and a theorem related to tensor ranks, which were also given in [17]. They are useful in deriving the results in Section 4. The first lemma indicates that the rank of a tensor can not be smaller than any of its $n$-ranks:

LEMMA 3.1. *Let $\mathcal{A} \in \mathbb{R}^{d_1 \times d_2 \times \cdots \times d_N}$ be an order-N tensor. Then*

$$\text{rank}_n(\mathcal{A}) \leq \min\{\text{rank}(\mathcal{A}), d_n\}, \quad \text{for } n = 1, 2, \ldots, N.$$

The next lemma illustrates a way to construct higher order tensors while preserving the rank.

LEMMA 3.2. *Let $\mathcal{A}$ be a tensor and $x$ be a non-zero vector. Then*

$$\text{rank}(\mathcal{A}) = \text{rank}(\mathcal{A} \otimes x).$$

The following lemma indicates that given any dimension $d_1 \times d_2 \times \cdots \times d_N$, we can construct a tensor of arbitrary rank $R \leq \min\{d_1, d_2, \ldots, d_N\}$.

LEMMA 3.3. *For $n = 1, \ldots, N$, let $x_1^{(n)}, \ldots, x_R^{(n)} \in \mathbb{R}^{d_n}$ be linearly independent. Then the tensor*

$$\mathcal{A} = \sum_{i=1}^{R} x_i^{(1)} \otimes x_i^{(2)} \otimes \cdots \otimes x_i^{(N)}$$

*has rank R.*

The next theorem is due to JáJá and Takche [23]. They showed that if $\mathcal{A}$ and $\mathcal{B}$ are order-3 tensors and at least one of them is a "stack" of two matrices, then the rank of their direct sum is equal to the sum of their ranks.

THEOREM 3.4 (JáJá–Takche). *Let* $\mathcal{A} \in \mathbb{R}^{d_1 \times d_2 \times d_3}$ *and* $\mathcal{B} \in \mathbb{R}^{c_1 \times c_2 \times c_3}$. *If* $2 \in \{d_1, d_2, d_3, c_1, c_2, c_3\}$, *then*

$$\mathrm{rank}(\mathcal{A} \oplus \mathcal{B}) = \mathrm{rank}(\mathcal{A}) + \mathrm{rank}(\mathcal{B}).$$

**3.1. The ill-posedness of the optimal low rank approximation problem.** de Silva and Lim [17] proved that for any order $N \geq 3$ and dimensions $d_1, \ldots, d_N \geq 2$, there exists a rank-$(r+1)$ tensor that has no optimal rank-$r$ approximation, for any $r = 2, \ldots, \min\{d_1, \ldots, d_N\}$. This result was further generalized to an arbitrary rank gap, i.e., there exists a rank-$(r+s)$ tensor that has no optimal rank-$r$ approximation, for some $r$'s and $s$'s.

Essentially, this ill-posedness of the optimal approximation problem is illustrated by the fact that the tensor

$$\mathcal{E} := e_2 \otimes e_1 \otimes e_1 + e_1 \otimes e_2 \otimes e_1 + e_1 \otimes e_1 \otimes e_2 \in \mathbb{R}^{2 \times 2 \times 2},$$

where $e_i$ is the $i$-th column of the identity matrix, has rank 3 but can be approximated arbitrarily closely by rank-at-most-2 tensors. Hence $\mathcal{E}$ does not have an optimal rank-2 approximation. Then according to Theorem 3.4 and Lemma 3.2, the ill-posedness of the problem can be generalized to arbitrary rank and order, by constructing higher rank and higher order tensors using direct sums and tensor products. We restate one of the results of [17] in the following theorem. For details of the proof, see the original paper.

THEOREM 3.5. *For* $N \geq 3$ *and* $d_1, d_2, \ldots, d_N \geq 2$, *there exists a tensor* $\mathcal{A} \in \mathbb{R}^{d_1 \times d_2 \times \cdots \times d_N}$ *of rank* $r + s$ *that has no optimal rank-r approximation, for any* $r$ *and* $s \geq 1$ *satisfying* $2s \leq r \leq \min\{d_1, d_2, \ldots, d_N\}$.

**4. The tensor SVD and its (non-)existence.** The definition used for the singular value decomposition of a tensor generalizes the matrix SVD from the angle of an expansion of outer product matrices, which becomes an expansion into a sum of rank-1 tensors.

DEFINITION 4.1. *If a tensor* $\mathcal{A} \in \mathbb{R}^{d_1 \times d_2 \times \cdots \times d_N}$ *can be written in the form*

$$\mathcal{A} = \sum_{i=1}^{R} \sigma_i u_i^{(1)} \otimes u_i^{(2)} \otimes \cdots \otimes u_i^{(N)}, \tag{4.1}$$

*where* $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_R > 0$ *and* $\left\langle u_j^{(n)}, u_k^{(n)} \right\rangle = \delta_{jk}$ *(Kronecker delta) for* $n = 1, 2, \ldots, N$, *then* (4.1) *is said to be the* tensor singular value decomposition (tensor SVD) *of* $\mathcal{A}$. *The* $\sigma_i$'s *are* singular values *and the* $u_i^{(n)}$'s *for* $i = 1, \ldots, R$ *are the* mode-$n$ singular vectors.

We also call (4.1) the *SVD* of tensor $\mathcal{A}$ for short where there is no ambiguity about tensors and matrices. In fact, when $N = 2$, i.e., $\mathcal{A}$ is a matrix, the tensor SVD of $\mathcal{A}$ boils down to the matrix SVD. Expression (4.1) can equivalently be written in the form

$$\mathcal{A} = \mathcal{D} \times_1 U^{(1)} \times_2 U^{(2)} \times \cdots \times_N U^{(N)}, \tag{4.2}$$

where $\mathcal{D} \in \mathbb{R}^{R \times R \times \cdots \times R}$ is the diagonal *core tensor* with $\mathcal{D}_{ii\ldots i} = \sigma_i$, and

$$U^{(n)} = \left[ u_1^{(n)}, u_2^{(n)}, \ldots, u_R^{(n)} \right] \in \mathbb{R}^{d_n \times R} \tag{4.3}$$

are orthogonal matrices for $n = 1, 2, \ldots, N$. The following proposition indicates that the tensor SVD is rank revealing.

PROPOSITION 4.2. *If a tensor $\mathcal{A}$ has the SVD as* (4.1)*, then* $\mathrm{rank}(\mathcal{A}) = R$.

*Proof.* This follows from Lemma 3.3. □

Trivially, if a tensor is constructed as in (4.1), its SVD exists. However, in general, a tensor of order $N \geq 3$ may fail to have such a decomposition. In this section, we identify some of these situations.

To begin with, note that the orthogonality of each $U^{(n)}$ implies that $R \leq d_n$ for each $n$, i.e., $R \leq \min\{d_1, d_2, \ldots, d_N\}$. This leads to the following simple result.

PROPOSITION 4.3. *A tensor $\mathcal{A} \in \mathbb{R}^{d_1 \times d_2 \times \cdots \times d_N}$ with* $\mathrm{rank}(\mathcal{A}) > \min\{d_1, d_2, \ldots, d_N\}$ *does not admit a tensor SVD.*

*Proof.* The existence of a tensor SVD such as in (4.1) would trivially lead to a contradiction since the tensor in (4.1) has rank $R$ with $R \leq \min\{d_1, d_2 \cdots, d_N\}$. □

Note that Theorem 3.5 guarantees that the condition of Proposition 4.3 is not vacuously satisfied, for any order $N \geq 3$ and dimensions $d_1, d_2, \ldots, d_N \geq 2$.

COROLLARY 4.4. *Given a tensor $\mathcal{A}$ satisfying the condition in Proposition 4.3, any tensor of the form*

$$\mathcal{A} \otimes x^{(N+1)} \otimes \cdots \otimes x^{(N+\ell)},$$

*where $\ell \geq 1$ and $x^{(N+1)}, \ldots, x^{(N+\ell)}$ are nonzero vectors, does not admit a tensor SVD.*

*Proof.* This follows from Proposition 4.3 and Lemma 3.2. □

COROLLARY 4.5. *A tensor $\mathcal{A} \in \mathbb{R}^{d_1 \times d_2 \times \cdots \times d_N}$ does not admit a tensor SVD if there exists at least one mode $n$ such that* $\mathrm{rank}_n(\mathcal{A}) > \min\{d_1, d_2, \ldots, d_N\}$.

*Proof.* This follows from Proposition 4.3 and Lemma 3.1. □

PROPOSITION 4.6. *There exists a tensor $\mathcal{A} \in \mathbb{R}^{d_1 \times d_2 \times \cdots \times d_N}$ which does not admit a tensor SVD whenever*

$$d := \max_n\{d_n\} > \min_n\{d_n\} \quad \textit{and} \quad d^2 \leq \prod_{n=1}^{N} d_n.$$

*Proof.* Without loss of generality, assume that $d = d_1 \geq d_2 \geq \cdots \geq d_N$ and let $d' = d_2 \times \cdots \times d_N$. Since $d \leq d'$, for an arbitrary set of orthonormal vectors $\{a_i \in \mathbb{R}^{d'} \mid i = 1, \ldots, d\}$, we can construct a tensor $\mathcal{A}$ whose unfolding $A_{(1)} = [a_1, a_2, \ldots, a_d]^T$. Then $\mathrm{rank}_1(\mathcal{A}) = d$. By Corollary 4.5, $\mathcal{A}$ does not admit a tensor SVD. □

Note that when $N = 2$, i.e., for the matrix case, it is impossible for $d_1$ and $d_2$ to satisfy the condition in the proposition.

In closing this section, we provide a necessary and sufficient condition to characterize the existence of the tensor SVD.[3] This is related to the HOSVD proposed by [12]. The essential relation underlying the theorem is that the mode-$n$ singular vectors of $\mathcal{A}$, when its SVD exists, are also the left singular vectors of the unfolding $A_{(n)}$.

---

[3] As pointed out by a referee, the provided relation may have long been recognized in other fields of research, such as signal processing, at least in the case of distinct singular values.

THEOREM 4.7. *A tensor $\mathcal{A}$ admits an SVD if and only if there exists an HOSVD of $\mathcal{A}$ such that the core is diagonal.*

*Proof.* The sufficient condition is obvious. Consider the necessary condition. If $\mathcal{A}$ can be written in the form (4.1), define the tensor

$$\mathcal{W}_i^{(n)} := u_i^{(n+1)} \otimes \cdots \otimes u_i^{(N)} \otimes u_i^{(1)} \otimes \cdots \otimes u_i^{(n-1)},$$

and let $w_i^{(n)}$ be the vectorization of $\mathcal{W}_i^{(n)}$. Then the unfolding of $\mathcal{A}$ along mode $n$ is

$$A_{(n)} = \sum_{i=1}^{R} \sigma_i u_i^{(n)} \otimes w_i^{(n)}.$$

Since $\left\langle u_j^{(n)}, u_k^{(n)} \right\rangle = \delta_{jk}$ for all $n$, we have $\left\langle w_j^{(n)}, w_k^{(n)} \right\rangle = \delta_{jk}$. Hence the above form is the SVD of matrix $A_{(n)}$. In other words, the vectors $u_1^{(n)}, \ldots, u_R^{(n)}$ are the left singular vectors of $A_{(n)}$. From the construction of the HOSVD, Equation (4.2) is a valid HOSVD for $\mathcal{A}$.[4] $\square$

The proof of the above theorem indicates that if the SVD of a tensor $\mathcal{A}$ exists, its singular values coincide with the nonzero mode-$n$ singular values in its HOSVD. However the HOSVD of a tensor may not be unique, since the SVD of the unfoldings $A_{(n)}$'s are not guaranteed to be unique. Hence even if a tensor is constructed as in (4.1), its HOSVD will not necessarily recover this form. This is the reason why in the above theorem we use the phase "... if there exists ...".

It is interesting to note again that the non-uniqueness of matrix SVD is caused by duplicate singular values, however the tensor SVD is unique (if it exists) even when some of the singular values are the same [43, Theorem 3.2].

**5. The optimal low rank orthogonal approximation.** The problem addressed by tensor analysis is to approximate some tensor $\mathcal{A}$ by a linear combination of tensors $\mathcal{T}_1, \mathcal{T}_2, \ldots, \mathcal{T}_r$ that have "special" structures, e.g., rank-1 tensors, orthogonal tensors, or simple tensors[5]. For this it is desirable to minimize

$$\left\| \mathcal{A} - \sum_{i=1}^{r} \sigma_i \mathcal{T}_i \right\|_F$$

for a given $r$. Without loss of generality, we assume that $\|\mathcal{T}_i\|_F = 1$ for all $i$. As discussed in Section 3.1, if the $\mathcal{T}_i$'s are rank-1 tensors, the infimum of the above expression might not necessarily be attained. The following proposition reveals some properties when the infimum is indeed achieved.

PROPOSITION 5.1. *Given a tensor $\mathcal{A}$ and a positive integer $r$, consider a set of linear combinations of tensors of the form*

$$\mathcal{T} := \sum_{i=1}^{r} \sigma_i \mathcal{T}_i \tag{5.1}$$

*where the $\mathcal{T}_i$'s are arbitrary unit tensors. If $\inf \|\mathcal{A} - \mathcal{T}\|_F$ is reached on this set, then for the optimal $\mathcal{T}$ and $\mathcal{T}_i$'s,*

$$\langle \mathcal{A} - \mathcal{T}, \mathcal{T}_i \rangle_F = 0 \quad \text{for } i = 1, 2, \ldots, r.$$

---

[4]In order to strictly conform to the definition of the HOSVD defined in [12], in (4.2) the size of $\mathcal{D}$ should be enlarged from $R \times R \times \cdots \times R$ to $d_1 \times d_2 \times \cdots \times d_N$ by padding zeros, and the $U^{(n)}$'s should be padded with orthogonal columns to make square shapes.

[5]A tensor is simple if it is the tensor product of two tensors.

*Furthermore, if the $\mathcal{T}_i$'s are required to be mutually F-orthogonal, then the optimal $\sigma_i$'s are related to the optimal $\mathcal{T}_i$'s by*

$$\sigma_i = \langle \mathcal{A}, \mathcal{T}_i \rangle_F \quad \text{for } i = 1, 2, \ldots, r. \tag{5.2}$$

*In this situation,*

$$\|\mathcal{T}\|_F = \sqrt{\sum_{i=1}^r \sigma_i^2}. \quad \text{and} \quad \|\mathcal{A} - \mathcal{T}\|_F^2 = \|\mathcal{A}\|_F^2 - \|\mathcal{T}\|_F^2. \tag{5.3}$$

*Proof.* If the infimum is attained by a certain set of $\sigma_i$'s and $\mathcal{T}_i$'s, and if there is a $j$ such that $\langle \mathcal{A} - \mathcal{T}, \mathcal{T}_j \rangle_F = \epsilon \neq 0$, then

$$\begin{aligned}
&\|\mathcal{A} - \textstyle\sum_{i=1}^r \sigma_i \mathcal{T}_i - \epsilon \mathcal{T}_j\|_F^2 \\
&= \|\mathcal{A} - \textstyle\sum_{i=1}^r \sigma_i \mathcal{T}_i\|_F^2 - 2\epsilon \langle \mathcal{A} - \textstyle\sum_{i=1}^r \sigma_i \mathcal{T}_i, \mathcal{T}_j \rangle_F + \epsilon^2 \|\mathcal{T}_j\|_F^2 \\
&= \|\mathcal{A} - \textstyle\sum_{i=1}^r \sigma_i \mathcal{T}_i\|_F^2 - \epsilon^2 < \|\mathcal{A} - \textstyle\sum_{i=1}^r \sigma_i \mathcal{T}_i\|_F^2,
\end{aligned}$$

which contradicts the assumption.

If the unit tensors $\mathcal{T}_i$'s are mutually F-orthogonal, then

$$0 = \langle \mathcal{A} - \textstyle\sum_{i=1}^r \sigma_i \mathcal{T}_i, \mathcal{T}_j \rangle_F = \langle \mathcal{A}, \mathcal{T}_j \rangle_F - \sigma_j \langle \mathcal{T}_j, \mathcal{T}_j \rangle_F = \langle \mathcal{A}, \mathcal{T}_j \rangle_F - \sigma_j.$$

Also,

$$\|\mathcal{T}\|_F^2 = \langle \mathcal{T}, \mathcal{T} \rangle_F = \sum_{i,j=1}^r \langle \sigma_i \mathcal{T}_i, \sigma_j \mathcal{T}_j \rangle_F = \sum_{i=1}^r \sigma_i^2,$$

and

$$\left\| \mathcal{A} - \sum_{i=1}^r \sigma_i \mathcal{T}_i \right\|_F^2 = \|\mathcal{A}\|_F^2 - \sum_{i=1}^r 2\sigma_i \langle \mathcal{A}, \mathcal{T}_i \rangle_F + \sum_{i=1}^r \sigma_i^2 = \|\mathcal{A}\|_F^2 - \sum_{i=1}^r \sigma_i^2 = \|\mathcal{A}\|_F^2 - \|\mathcal{T}\|_F^2.$$

□

The last part of the proof indicates that the equalities in (5.3) follow from the orthogonality of the $\mathcal{T}_i$'s and the relations (5.2). They do not require optimality.

In this section, we will see that if the $\mathcal{T}_i$'s are mutually orthogonal rank-1 tensors, then the infimum in the proposition can be attained. Formally, we will prove that the problem

$$\begin{aligned}
\min \quad & E := \left\| \mathcal{A} - \sum_{i=1}^r \sigma_i u_i^{(1)} \otimes u_i^{(2)} \otimes \cdots \otimes u_i^{(N)} \right\|_F \\
\text{s.t.} \quad & \left\langle u_j^{(n)}, u_k^{(n)} \right\rangle = \delta_{jk}, \quad \text{for } n = 1, 2, \ldots, N,
\end{aligned} \tag{5.4}$$

has a solution for any $\mathcal{A} \in \mathbb{R}^{d_1 \times d_2 \times \cdots \times d_N}$ and any $r \leq \min\{d_1, d_2, \ldots, d_N\}$. The solution for the case $r = \min\{d_1, d_2, \ldots, d_N\}$ leads to a decomposition of $\mathcal{A}$ where the diagonal of the core is maximized.

**5.1. Existence of the global optimum.** Let

$$\mathcal{T}_i := u_i^{(1)} \otimes u_i^{(2)} \otimes \cdots \otimes u_i^{(N)}, \quad \text{for } i = 1, \ldots, r, \tag{5.5}$$

and $\sigma_i$'s be defined as in (5.2), then according to Proposition 5.1 (see comments following the proof),

$$E^2 = \left\| \mathcal{A} - \sum_{i=1}^{r} \sigma_i \mathcal{T}_i \right\|_F^2 = \|\mathcal{A}\|_F^2 - \sum_{i=1}^{r} \sigma_i^2.$$

Hence minimizing $E$ is equivalent to maximizing $\sum_{i=1}^{r} \sigma_i^2$, i.e., the optimization problem (5.4) is equivalent to the following:

$$\max \quad E' := \sum_{i=1}^{r} \left( \mathcal{A} \times_1 u_i^{(1)^T} \times_2 u_i^{(2)^T} \times \cdots \times_N u_i^{(N)^T} \right)^2 \tag{5.6}$$

$$\text{s.t.} \quad \left\langle u_j^{(n)}, u_k^{(n)} \right\rangle = \delta_{jk}, \quad \text{for } n = 1, 2, \ldots, N.$$

Let

$$U^{(n)} = \left[ u_1^{(n)}, u_2^{(n)}, \ldots, u_r^{(n)} \right] \in \Omega^{(n)} \tag{5.7}$$

where

$$\Omega^{(n)} := \{ W \in \mathbb{R}^{d_n \times r} \mid W^T W = I \} \tag{5.8}$$

for $n = 1, 2, \ldots, N$. The problem (5.6) can be interpreted as that of maximizing $E'$ within the feasible region

$$\Omega := \Omega^{(1)} \times \Omega^{(2)} \times \cdots \times \Omega^{(N)}. \tag{5.9}$$

Since for each $n$ the set $\Omega^{(n)}$ is compact (see, e.g., [22, p. 69]), by Tychonoff Theorem, the feasible region $\Omega$ is compact. Under the continuous mapping $E'$, the image $E'(\Omega)$ is also compact. Hence it has a maximum. This proves the following theorem:

THEOREM 5.2. *There exists a solution to the problem (5.6) (or equivalently (5.4) with $\sigma_i$ defined in (5.2)) for any $r \leq \min\{d_1, d_2, \ldots, d_N\}$.*

**5.2. Relation to tensor decompositions.** Let $U^{(n)}$, $n = 1, \ldots, N$ be the solution to the problem (5.4) with $r = \min\{d_1, d_2, \ldots, d_N\}$ and $\sigma_i$ be defined in (5.2). Also for $n = 1, \ldots, N$, let $U^{(n)\perp}$ be a $d_n \times (d_n - r)$ matrix such that the square matrix

$$\tilde{U}^{(n)} := \left[ U^{(n)}, U^{(n)\perp} \right] \in \mathbb{R}^{d_n \times d_n} \tag{5.10}$$

is orthogonal. Further, define the tensor

$$\mathcal{S} := \mathcal{A} \times_1 \tilde{U}^{(1)^T} \times_2 \tilde{U}^{(2)^T} \times \cdots \times_N \tilde{U}^{(N)^T} \in \mathbb{R}^{d_1 \times d_2 \times \cdots \times d_N}. \tag{5.11}$$

Then the equality

$$\mathcal{A} = \mathcal{S} \times_1 \tilde{U}^{(1)} \times_2 \tilde{U}^{(2)} \times \cdots \times_N \tilde{U}^{(N)} \tag{5.12}$$

holds. This decomposition of $\mathcal{A}$ has the following two properties:

(i) The side-matrices $\tilde{U}^{(n)}$'s are orthogonal.

(ii) The (squared) norm of the diagonal of the core $\mathcal{S}$:

$$\sum_{i=1}^{\min\{d_1,\ldots,d_N\}} s_{ii\ldots i}^2 = \sum_{i=1}^{r} \left( \mathcal{A} \times_1 u_i^{(1)^T} \times_2 u_i^{(2)^T} \times \cdots \times_N u_i^{(N)^T} \right)^2 = \sum_{i=1}^{r} \sigma_i^2$$

is maximized among all the choices of the orthogonal side-matrices. This is known as *maximal diagonality* in [12].

**5.3. First order condition.** The Lagrangian of (5.6) is

$$L = \sum_{i=1}^{r} \sigma_i^2 - \sum_{j,k=1}^{r} \sum_{n=1}^{N} \mu_{j,k}^n \left( \left\langle u_j^{(n)}, u_k^{(n)} \right\rangle - \delta_{jk} \right), \tag{5.13}$$

where

$$\sigma_i = \mathcal{A} \times_1 u_i^{(1)^T} \times_2 u_i^{(2)^T} \times \cdots \times_N u_i^{(N)^T} \tag{5.14}$$

and the $\mu_{j,k}^n$'s are Lagrange multipliers. Define the vector

$$\begin{aligned} v_i^{(n)} &:= \mathcal{A} \times_1 u_i^{(1)^T} \times \cdots \times_{n-1} u_i^{(n-1)^T} \times_{n+1} u_i^{(n+1)^T} \times \cdots \times_N u_i^{(N)^T} \\ &\in \mathbb{R}^{1 \times \cdots \times 1 \times d_n \times 1 \cdots \times 1}. \end{aligned} \tag{5.15}$$

(Here we abuse the use of notation "=". More precisely, $v_i^{(n)}$ should be the mode-$n$ unfolding of the tensor on the right-hand side of (5.15).) It is not hard to see that

$$\left\langle u_i^{(n)}, v_i^{(n)} \right\rangle = \sigma_i$$

for all $n$ and $i$, and $v_i^{(n)}$ is the partial derivative of $\sigma_i$ with respect to $u_i^{(n)}$.

The partial derivative of the Lagrangian with respect to $u_i^{(n)}$ is

$$\frac{\partial L}{\partial u_i^{(n)}} = 2\sigma_i v_i^{(n)} - \sum_{j=1}^{r} \mu_{j,i}^n u_j^{(n)} - \sum_{k=1}^{r} \mu_{i,k}^n u_k^{(n)},$$

for any $n$ and $i$. By setting the partial derivatives to 0 and putting all equations related to the same $n$ in matrix form, we obtain the following equations:

$$\begin{bmatrix} v_1^{(n)} & \cdots & v_r^{(n)} \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{bmatrix} = \begin{bmatrix} u_1^{(n)} & \cdots & u_r^{(n)} \end{bmatrix} \begin{bmatrix} \frac{\mu_{1,1}^n + \mu_{1,1}^n}{2} & \cdots & \frac{\mu_{1,r}^n + \mu_{r,1}^n}{2} \\ \vdots & \ddots & \vdots \\ \frac{\mu_{r,1}^n + \mu_{1,r}^n}{2} & \cdots & \frac{\mu_{r,r}^n + \mu_{r,r}^n}{2} \end{bmatrix}, \tag{5.16}$$

for all $n = 1, 2, \ldots, N$. Let

$$V^{(n)} := \begin{bmatrix} v_1^{(n)}, v_2^{(n)}, \ldots, v_r^{(n)} \end{bmatrix}, \tag{5.17}$$

$$\Sigma := \operatorname{diag}(\sigma_1, \ldots, \sigma_r), \tag{5.18}$$

and let $M^{(n)}$ be the second term on the right-hand side of (5.16). Then (5.16) is compactly represented as

$$V^{(n)} \Sigma = U^{(n)} M^{(n)}, \quad n = 1, 2, \ldots, N. \tag{5.19}$$

In summary, the necessary condition of an extremum of the Lagrangian is the equation (5.19), where $V^{(n)}$ is defined in (5.17), $\Sigma$ is defined in (5.18), $U^{(n)}$ is defined in (5.7), and $M^{(n)}$ is symmetric, for all $n = 1, 2, \ldots, N$.

**5.4. Algorithm: LROAT.** We seek orthogonal matrices $U^{(n)}$'s and symmetric matrices $M^{(n)}$'s which satisfy the system (5.19). (The $\Sigma$ and $V^{(n)}$'s are computed from the $U^{(n)}$'s.) Note that the pair $U^{(n)}, M^{(n)}$ happens to be the polar decomposition of the matrix $V^{(n)}\Sigma$. Hence the system can be solved in an iterative fashion: We begin with an initial guess of the set of orthogonal matrices $\{U^{(1)}, U^{(2)}, \ldots, U^{(N)}\}$, which can be obtained, say, by the HOSVD of $\mathcal{A}$. For each $n$, we compute $V^{(n)}$ and $\Sigma$, and update $U^{(n)}$ as an orthogonal polar factor of $V^{(n)}\Sigma$. This procedure is iterated until convergence is observed. Algorithm 1 (LROAT) summarizes this idea.

---

**Algorithm 1** Low Rank Orthogonal Approximation of Tensors (LROAT)

---

**Input:** Tensor $\mathcal{A}$, rank $r$, orthogonal matrices $U^{(1)}, \ldots, U^{(N)}$ as initial guess
**Output:** $\sigma_1, \ldots, \sigma_r, U^{(1)}, \ldots, U^{(N)}$

 1: **repeat**
 2:    **for** $n \leftarrow 1, \ldots, N$ **do**
 3:       Compute $V^{(n)} = \left[v_1^{(n)}, \ldots, v_r^{(n)}\right]$ according to (5.15)
 4:       Compute $\Sigma = \text{diag}(\sigma_1, \ldots, \sigma_r)$ according to (5.14)
 5:       $[Q^{(n)}, H^{(n)}] \leftarrow$ polar-decomp$(V^{(n)}\Sigma)$
 6:       Update $U^{(n)} \leftarrow Q^{(n)}$
 7:    **end for**
 8: **until** convergence

---

Note that when $r = 1$, the matrix $V^{(n)} = \left[v_1^{(n)}\right]$ and $U^{(n)} = \left[u_1^{(n)}\right]$, which means that for each iteration $u_1^{(n)}$ is updated as the normalized $v_1^{(n)}$. This indicates that the LROAT algorithm for $r = 1$ boils down to the ALS method [43] (or the so-called higher-order power method [13, 27]) for computing the optimal rank-1 approximation. Hence, it is not unexpected to see in the numerical experiments that in general LROAT converges linearly. We also comment that LROAT is not an alternating least squares method (except for the case $r = 1$) by the nature of the update of $U^{(n)}$.

**5.5. Convergence analysis.** LROAT employs an alternating procedure (iterating through $U^{(1)}, U^{(2)}, \ldots, U^{(N)}$), where in each step all but one ($U^{(n)}$) parameters are fixed. In general, algorithms of this type, including alternating least squares, are not guaranteed to converge. Specifically, the objective function may converge but not the parameters. (See, for example, [32] for some discussions.) For LROAT, we are also unable yet to prove the global convergence, though empirically it appears to hold. However, in this section, we will prove that: (1) The iterations monotonically increase the objective value $E'$ (Theorem 5.4); (2) Under a mild condition, of the generated parameter sequence, every converging subsequence converges to a stationary point of the objective function (Theorem 5.7); and (3) In a neighborhood of a local maximum, the parameter sequence converges to this stationary point (Theorem 5.9).

Before analyzing the convergence behavior of LROAT, we index all the iterates. The outer-loop is indexed by $p$ and the overall iterations are indexed by $idx$, which is equal to $n + (p - 1)N$. In other words, Algorithm 1 is rewritten as follows. In particular, the numbered lines correspond to the lines in Algorithm 1.

---

   **for** $p \leftarrow 1, 2, \ldots$ **do**
      **for** $n \leftarrow 1, \ldots, N$ **do**
        $idx = n + (p - 1)N$

For all $i$, compute $\sigma_i^{(idx)}$ according to $U_{(p+1)}^{(1)}, \ldots, U_{(p+1)}^{(n-1)}, U_{(p)}^{(n)}, U_{(p)}^{(n+1)}, \ldots, U_{(p)}^{(N)}$

Objective $E'^{(idx)} = \sum_{i=1}^{r} \left( \sigma_i^{(idx)} \right)^2$

3:    Compute $V_{(p)}^{(n)}$ from $U_{(p+1)}^{(1)}, \ldots, U_{(p+1)}^{(n-1)}, U_{(p)}^{(n+1)}, \ldots, U_{(p)}^{(N)}$

4:    Assign $\Sigma^{(idx)} = \text{diag}\left( \sigma_1^{(idx)}, \ldots, \sigma_r^{(idx)} \right)$

5:    Polar decomposition $V_{(p)}^{(n)} \Sigma^{(idx)} = Q_{(p)}^{(n)} H_{(p)}^{(n)}$

6:    Update $U_{(p+1)}^{(n)} = Q_{(p)}^{(n)}$

   **end for**

 **end for**

---

The following lemma, which is well-known when the matrix $A$ is square, reveals the trace maximizing property that is important for the convergence analysis of LROAT.

LEMMA 5.3. *Let matrix $A \in \mathbb{R}^{m \times n}$, $m \geq n$, have the polar decomposition $A = QH$ where $Q \in \mathbb{R}^{m \times n}$ is the orthogonal polar factor and $H \in \mathbb{R}^{n \times n}$ is the symmetric positive semi-definite factor, then*

$$\max_{P \in \mathbb{R}^{m \times n}, P^T P = I} \text{tr}(P^T A)$$

*is attained when $P = Q$.*

*Proof.* Any $P$ can be written as $ZQ$, where $Z \in \mathbb{R}^{m \times m}$ is orthogonal. Then

$$\text{tr}(P^T A) = \text{tr}(Q^T Z^T Q H) = \text{tr}(Z^T Q H Q^T).$$

Since $QHQ^T$ is symmetric positive semi-definite, $\max \text{tr}(Z^T Q H Q^T)$ is attained when $Z = I$. $\square$

Since $U_{(p+1)}^{(n)}$ is the orthogonal polar factor of $V_{(p)}^{(n)} \Sigma^{(idx)}$, by Lemma 5.3,

$$\sum_{i=1}^{r} \left( \sigma_i^{(idx)} \right)^2 = \text{tr}\left( U_{(p)}^{(n)T} V_{(p)}^{(n)} \Sigma^{(idx)} \right) \leq \text{tr}\left( U_{(p+1)}^{(n)}{}^T V_{(p)}^{(n)} \Sigma^{(idx)} \right) = \sum_{i=1}^{r} \sigma_i^{(idx+1)} \sigma_i^{(idx)}.$$

Then by the Cauchy-Schwarz inequality,

$$\sum_{i=1}^{r} \left( \sigma_i^{(idx)} \right)^2 \leq \sum_{i=1}^{r} \sigma_i^{(idx+1)} \sigma_i^{(idx)} \leq \sum_{i=1}^{r} \left( \sigma_i^{(idx+1)} \right)^2, \tag{5.20}$$

and

$$\sum_{i=1}^{r} \left( \sigma_i^{(idx)} \right)^2 = \sum_{i=1}^{r} \left( \sigma_i^{(idx+1)} \right)^2 \quad \text{iff} \quad \sigma_i^{(idx)} = \sigma_i^{(idx+1)} \text{ for all } i. \tag{5.21}$$

Inequality (5.20) means that each update of $U^{(n)}$ increases the value of the objective function $E'$, i.e.,

$$E'^{(idx)} \leq E'^{(idx+1)}.$$

Since $E'$ is bounded from above (existence of the maximum, see Theorem 5.2), the sequence $\{E'^{(idx)}\}_{idx=1}^{\infty}$ converges. Note that the convergence does not depend on the initial guess input to the algorithm. Formally, we have established the following theorem:

THEOREM 5.4. *Given any initial guess, the iterations of Algorithm 1 monotonically increase the objective function $E'$ defined in (5.6) to a limit.*

The convergence of the objective function does not necessarily imply that the function parameters will converge. However, in our case since the parameters $U^{(n)}$'s are bounded, they admit converging subsequences. Next we will show that every such subsequence converges to a stationary point of $E'$. For this, the following lemma uses a helper function $f$.

LEMMA 5.5. *Let $T : \Theta \to \Theta$ be a continuous mapping and a sequence $\{\theta_n \in \Theta\}_{n=1}^{\infty}$ be generated from the fixed point iteration $\theta_{n+1} = T(\theta_n)$. If there exists a continuous function $f : \Theta \to \mathbb{R}$ satisfying the following two conditions:*

(i) *The sequence $\{f(\theta_n)\}_{n=1}^{\infty}$ converges;*

(ii) *For $\theta \in \Theta$, if $f(T(\theta)) = f(\theta)$ then $T(\theta) = \theta$;*

*then every converging subsequence of $\{\theta_n\}_{n=1}^{\infty}$ converges to a fixed point of $T$.*

*Proof.* Let $\{\theta_{s_n}\}_{n=1}^{\infty}$ be a converging subsequence of $\{\theta_n\}_{n=1}^{\infty}$, where $\theta_{s_n} \to \theta^*$. Also let $f^*$ be the limit of $f(\theta_n)$. Then $f(\theta_{s_n}) \to f(\theta^*)$, therefore $f(\theta^*) = f^*$. Meanwhile from the continuity of $T$ and $f$, we have $T(\theta_{s_n}) \to T(\theta^*)$ and $f(\theta_{s_n+1}) = f(T(\theta_{s_n})) \to f(T(\theta^*))$, which implies that $f(T(\theta^*)) = f^*$. Condition (ii) of the lemma now implies that $\theta^* = T(\theta^*)$. ☐

Our objective function $E'$ is just one such helper function $f$, and the orthogonal polar factor function plays the role of the mapping $T$ in the above lemma. The following lemma establishes the fact that the orthogonal polar factor function is continuous.

LEMMA 5.6. *The orthogonal polar factor function $g : A \to Q$ defined on the set of matrices with full column rank is continuous. Here $Q$ is the orthogonal polar factor of $A \in \mathbb{R}^{m \times n}$, $m \geq n$.*

*Proof.* First, function $g$ is well defined, since the orthogonal polar factor of a full rank matrix exists and is unique [21]. If $Q$ and $Q'$ are the orthogonal polar factors of $A$ and $A'$, respectively, Sun and Chen [39] have shown that

$$\|Q - Q'\|_F \leq \frac{2}{\|A^+\|_2} \|A - A'\|_F,$$

where $^+$ means pseudo inverse. Hence if $A_1, A_2, \ldots$ converges to $A^*$, then $g(A_1), g(A_2), \ldots$ converges to $g(A^*)$. ☐

Now we are ready to prove the following result.

THEOREM 5.7. *Every converging subsequence of $\left\{ U_{(p)}^{(1)}, \ldots, U_{(p)}^{(N)} \right\}_{p=1}^{\infty}$ generated by Algorithm 1 converges to a stationary point of the objective function $E'$ defined in (5.6), provided the matrices $V^{(n)}$ in line 3 of the algorithm do not become rank-deficient throughout the iterations.*

*Proof.* For convenience, let $U$ denote the side-by-side concatenation of the $U^{(n)}$'s, i.e., at iteration number $p$ we write $U_{(p)} := \left[ U_{(p)}^{(1)}, \ldots, U_{(p)}^{(N)} \right]$. For each iteration, $V_{(p)}^{(n)} \Sigma^{(idx)}$ is computed from $U_{(p)}$ and polar factorized, and $U_{(p)}^{(n)}$ is updated. Let $T$ be the composite of all these iterations running $n$ from 1 to $N$. That is, $U_{(p+1)} = T(U_{(p)})$. It is not hard to see that $T$ is continuous by Lemma 5.6. The objective function $E'$ taking parameter $U_{(p)}$ has been previously shown such that the sequence $\{E'(U_{(p)})\}_{p=1}^{\infty}$ is monotonically converging.

Hence by Lemma 5.5, in order to prove this theorem, it will suffice to show that $E'(T(U)) = E'(U)$ implies $T(U) = U$. Then every converging subsequence of $\{E'(U_{(p)})\}_{p=1}^{\infty}$ converges to a fixed point, which satisfies the first order condition (5.19), i.e., it is also a stationary point of $E'$.

If $E'(T(U)) = E'(U)$, formula (5.21) indicates that the $\sigma_i$ values have not changed after the iteration. In particular, for any $n$, the update of $U^{(n)}$ has not changed $\text{tr}\left(U^{(n)T}V^{(n)}\Sigma\right)$. Since the orthogonal polar factor of $V^{(n)}\Sigma$ is unique when $V^{(n)}$ is not rank-deficient, this means that $U^{(n)}$ has not changed. This in turn means that $U$ is a fixed point of the mapping $T$. □

The condition in the theorem is not a strong requirement in general. Of course, the columns $v_i^{(n)}$ of the matrix $V^{(n)}$, as computed from (5.15), will be linearly dependent if the $n$-rank of $\mathcal{A}$ is less than $r$. For practical applications, the tensor usually has full $n$-ranks for all $n$, so this does not hamper the applicability of the theorem.

Though the global convergence of $\{U_{(p)}\}$ is not determined, when localized, it is possible that this parameter sequence converges. The following lemma and theorems consider this situation.

LEMMA 5.8. *If a sequence $\{\theta_n\}_{n=1}^\infty$ is bounded, and all of its converging subsequences converge to $\theta^*$, then $\theta_n \to \theta^*$.*

*Proof.* (By contradiction.) If $\{\theta_n\}_{n=1}^\infty$ does not converge to $\theta^*$, then there is an $\epsilon > 0$ such that there exists a subsequence $S = \{\theta_{s_n}\}_{n=1}^\infty$, where $\|\theta_{s_n} - \theta^*\| \geq \epsilon$ for all $n$. Since $S$ is bounded, it has a converging subsequence $S'$. Then $S'$ as a subsequence of $\{\theta_n\}_{n=1}^\infty$ converges to a limit other than $\theta^*$. □

THEOREM 5.9. *Let $U_* = \left[U_*^{(1)}, \ldots, U_*^{(N)}\right]$ be a local maximum of the objective function $E'$ defined in (5.6). If the sequence $\left\{U_{(p)} := \left[U_{(p)}^{(1)}, \ldots, U_{(p)}^{(N)}\right]\right\}_{p=1}^\infty$ generated by Algorithm 1 lies in a neighborhood of $U_*$, where $U_*$ is the only stationary point in that neighborhood, and if the full rank requirement in Theorem 5.7 is satisfied, then the sequence $\{U_{(p)}\}_{p=1}^\infty$ converges to $U_*$.*

*Proof.* This immediately follows from Theorem 5.7 and Lemma 5.8. □

Note that since the starting elements of a sequence have no effect on its convergence behavior, the above theorem holds whenever the tailing subsequence, starting from a sufficiently large $p$, lies within the neighborhood.

A weaker, but simpler, result is the following corollary.

COROLLARY 5.10. *Let $U_* = \left[U_*^{(1)}, \ldots, U_*^{(N)}\right]$ be a local maximum of the objective function $E'$ defined in (5.6). If this local maximum is unique and if the full rank requirement in Theorem 5.7 is satisfied, then the sequence $\left\{U_{(p)} := \left[U_{(p)}^{(1)}, \ldots, U_{(p)}^{(N)}\right]\right\}_{p=1}^\infty$ generated by Algorithm 1 converges to $U_*$.*

**5.6. LROAT for symmetric tensors.** An order-$N$ tensor $\mathcal{A} \in \mathbb{R}^{d \times d \times \cdots \times d}$, whose dimensions of all modes are the same, is *symmetric* if for all permutations $\pi$,

$$a_{i_1, i_2, \ldots, i_N} = a_{i_{\pi(1)}, i_{\pi(2)}, \ldots, i_{\pi(N)}}.$$

For symmetric tensors, usually the approximation problem (5.4) has an additional constraint that the side-matrices $U^{(n)}$'s are the same for all $n$, i.e., the problem becomes

$$\min \quad E = \left\|\mathcal{A} - \sum_{i=1}^r \sigma_i u_i \otimes u_i \otimes \cdots \otimes u_i\right\|_F \tag{5.22}$$

$$\text{s.t.} \quad \langle u_j, u_k \rangle = \delta_{jk}.$$

Applying similar arguments to those in Section 5.1, it is easily seen that (5.22) is

equivalent to the following problem:

$$\max \quad E' = \sum_{i=1}^{r} \left( \mathcal{A} \times_1 u_i^T \times_2 u_i^T \times \cdots \times_N u_i^T \right)^2 \tag{5.23}$$

$$\text{s.t.} \quad \langle u_j, u_k \rangle = \delta_{jk}.$$

The supremum of $E'$ can be attained. Further, the "maximal-diagonality" decomposition of $\mathcal{A}$ (c.f. Equation (5.12)) has an additional property that the core $\mathcal{S}$ is symmetric. Also, the first order condition (5.19) is simplified to

$$V\Sigma = UM.$$

Hence, there are two approaches to compute the approximation for the symmetric tensor $\mathcal{A}$. The first approach is to directly apply LROAT on $\mathcal{A}$. Theorems in the above section guarantee the convergence under mild assumptions, but the side-matrices might no longer be the same, though in the next section an experiment indicates that they indeed converge to the same matrix. The second approach is to only use a single initial guess $U$ and omit the for-loop on $n$ (line 2 of Algorithm 1). We call this the *symmetric variant of LROAT*. In this case Theorem 5.4 no longer holds, i.e., the iterations might not monotonically increase the objective value $E'$ defined in (5.23), since the for-loop on $n$ is omitted. An experiment in the next section shows an oscillating phenomenon, which is similar to the one indicated in Figure 4.1 of [27], for the objective value $E'$.

**6. Numerical experiments.** This section will show a few experiments to illustrate the convergence behavior and the approximation quality of LROAT. For comparisons are the ALS methods for Tucker and PARAFAC, whose implementations are based on the codes from the MATLAB Tensor Toolbox developed by Bader and Kolda [2]. We use the major left singular vectors of the unfoldings as the initial guess input for all the algorithms compared. When it comes to the quality of the final approximation, experience shows that compared with random orthonormal vectors, singular vectors as initial guesses do not offer any advantage. It has been argued that running the algorithms several times using different sets of random initial guess enhances the probability of hitting the global optimum. We use singular vectors here only for repeatability.

**6.1. Convergence of LROAT.** In the first experiment, we test LROAT (and the symmetric variant of LROAT mentioned in Section 5.6) on a few tensors listed in Table 6.1. The results are shown in Figures 6.1 and 6.2. Each row of the figures is one test on a tensor. The left plot shows the objective value $E'$ (the same as the norm of the approximated tensor $\mathcal{T}$) for each iteration $p$, while the right plot shows the convergence history of the $U^{(n)}$'s. Since the optima are unknown, we plot the values $\left\| U_{(p)}^{(n)} - U_{(p-1)}^{(n)} \right\|_F$ to indicate the convergence of the sequence $\left\{ U_{(p)}^{(n)} \right\}_{p=1,2,\ldots}$. Since these values are plotted on logarithimic scale, if the curves are bounded from above by a straight decreasing line, then it is indicated that the convergence of the sequence is at least linear.

Figures 6.1 and 6.2 show in a total five tests. The first test (Figure 6.1(a)) uses a randomly generated tensor $\mathcal{A}_1$. The second test (Figure 6.1(b)) uses a low-rank-plus-Gaussian-noise tensor

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

*The tensors used for the first experiment. The value $r$ is the rank input to LROAT; it is not the rank of the tensor.*

| Tensor | Dimensions | $r$ | Notes |
|--------|-----------|-----|-------|
| $\mathcal{A}_1$ | $20 \times 16 \times 10 \times 32$ | 5 | random tensor |
| $\mathcal{A}_2$ | $20 \times 16 \times 10 \times 32$ | 5 | rank-5 tensor + Gaussian noise |
| $\mathcal{A}_3$ | $10 \times 10 \times 10$ | 5 | the $(i, j, k)$-entry $= 1/(i^2 + j^2 + k^2)$ |
| $\mathcal{A}_4$ | $3 \times 3 \times 3 \times 3$ | 2 | see [27, Example 1] |

where the low rank tensor $\mathcal{B}_1$ is in the form (2.3) with $r = 5$, the Gaussian noise tensor $\mathcal{B}_2$ has normally distributed elements, and $\rho = 0.1 \, \|\mathcal{B}_1\|_F / \|\mathcal{B}_2\|_F$. In these two tests the two tensors are applied to the LROAT algorithm. The third test (Figure 6.1(c)) uses a symmetric tensor $\mathcal{A}_3$ with entries

$$(\mathcal{A}_3)_{ijk} = \frac{1}{i^2 + j^2 + k^2}.$$

In this test $\mathcal{A}_3$ is applied to the symmetric variant of LROAT. All the three tests show a linear convergence rate. The fourth (Figure 6.2(a)) and the fifth (Figure 6.2(b)) test use a symmetric tensor $\mathcal{A}_4$ introduced in [27, Example 1]:

$$
\begin{array}{lll}
(\mathcal{A}_4)_{1111} = 0.2883, & (\mathcal{A}_4)_{1112} = -0.0031, & (\mathcal{A}_4)_{1113} = 0.1973, \\
(\mathcal{A}_4)_{1112} = -0.2485, & (\mathcal{A}_4)_{1123} = -0.2939, & (\mathcal{A}_4)_{1133} = 0.3847, \\
(\mathcal{A}_4)_{1222} = 0.2972, & (\mathcal{A}_4)_{1223} = 0.1862, & (\mathcal{A}_4)_{1233} = 0.0919, \\
(\mathcal{A}_4)_{1333} = -0.3619, & (\mathcal{A}_4)_{2222} = 0.1241, & (\mathcal{A}_4)_{2223} = -0.3420, \\
(\mathcal{A}_4)_{2233} = 0.2127, & (\mathcal{A}_4)_{2333} = 0.2727, & (\mathcal{A}_4)_{3333} = -0.3054.
\end{array}
$$

In [27], the symmetric higher-order power method for computing the optimal rank-1 approximation of $\mathcal{A}_4$ is shown to be non-converging. We experiment with this tensor with $r = 2$ on LROAT and the symmetric variant of LROAT. Figure 6.2(a) shows that when applied to LROAT, the approximation to $\mathcal{A}_4$ indeed linearly converges, and what's more, all the side-matrices converge to the same result. The approximation computed by LROAT is

$$\mathcal{A}_4 \approx \sigma_1 u^{(1)} \otimes u^{(1)} \otimes u^{(1)} \otimes u^{(1)} + \sigma_2 u^{(2)} \otimes u^{(2)} \otimes u^{(2)} \otimes u^{(2)}$$

with

$$
\begin{array}{ll}
\sigma_1 = -1.0939, & u^{(1)} = \begin{bmatrix} -0.5946 & 0.7503 & 0.2890 \end{bmatrix}^T, \\
\sigma_2 = -0.55594, & u^{(2)} = \begin{bmatrix} 0.1947 & -0.2144 & 0.9572 \end{bmatrix}^T.
\end{array}
$$

On the other hand, Figure 6.2(b) shows that the symmetric variant of LROAT fails to converge.

**6.2. Low rank orthogonal approximation compared with Tucker and PARAFAC.** In the second experiment, we compare the approximation quality of three different models: Low rank orthogonal approximation (without confusion in this section, we call this model "LROAT", which happens to be the name of the algorithm, for short), Tucker and PARAFAC. See Figure 6.3. We experiment with two tensors: a low-rank-plus-Gaussian-noise tensor which is generated the same way

as $\mathcal{A}_2$ and a real-life tensor. The latter is obtained from a problem in acoustics [20], and the data can be downloaded from [16]. The residual norms

$$res(p) := \frac{\left\| \mathcal{A} - \mathcal{T}_{(p)} \right\|_F}{\left\| \mathcal{A} \right\|_F}$$

over all the iterations $p$ are plotted.

Figure 6.3 indicates three facts: (1) The three models approximate the data tensor well to some extent (less than 35% of the information is lost due to approximation); (2) PARAFAC is usually slow to converge; (3) The residual norm for LROAT is larger than those of Tucker and PARAFAC. The last fact is not unexpected since LROAT can be considered a special case of Tucker and of PARAFAC: The Tucker model has a full core while the core for LROAT is diagonal, and unlike LROAT the side-matrices in the PARAFAC model are not restricted to be orthogonal.

**6.3. An application.** In the blind source separation (BSS) problem [5], the cumulant tensor of order 4 is a rank-$R$ tensor:

$$\sum_{i=1}^{R} \sigma_i u_i \otimes u_i \otimes u_i \otimes u_i, \tag{6.1}$$

where $R$ is the number of sources and $u_i$ is the $i$-th column of the mixing matrix. In the prewhitening approach for the BSS problem, the $u_i$'s become the columns of the composite of the whitening matrix and the mixing matrix, that is, the $u_i$'s are length-$R$ vectors and are orthonormal. Hence, this prewhitening approach reduces to computing the tensor SVD of the cumulant tensor. Since in practice this tensor is estimated from a finite data set, it is not exact. Thus, the low rank orthogonal approximation becomes a suitable tool to recover the $u_i$'s.

In an experiment, we let $R = 3$ and generate a data tensor

$$\mathcal{A}_5 = \mathcal{B}_3 + \rho \mathcal{B}_4$$

where $\mathcal{B}_3$ is as (6.1), $\mathcal{B}_4$ is a symmetric tensor with normally distributed elements, and $\rho = 0.05 \left\| \mathcal{B}_3 \right\|_F / \left\| \mathcal{B}_4 \right\|_F$. The $\sigma_i$'s are

$$\sigma_1 = 0.7942, \quad \sigma_2 = 0.5678, \quad \sigma_3 = 0.4611,$$

and the $u_i$'s are

$$U = [u_1, u_2, u_3] = \begin{bmatrix} 0.0974 & 0.4049 & 0.9092 \\ 0.9918 & -0.1154 & -0.0548 \\ 0.0827 & 0.9071 & -0.4128 \end{bmatrix}.$$

We use four methods to compute the rank-$R$ (or rank-$(R, R, R)$) approximations to $\mathcal{A}_5$: LROAT, incremental rank-1 approximation, PARAFAC, and Tucker. All the four methods return same side-matrices for all modes. They are:

$$U_{\text{LROAT}} = \begin{bmatrix} 0.0937 & 0.3822 & 0.9193 \\ 0.9918 & -0.1164 & -0.0527 \\ 0.0869 & 0.9167 & -0.3899 \end{bmatrix}, \quad U_{\text{inc}} = \begin{bmatrix} 0.0841 & 0.3795 & 0.9162 \\ 0.9929 & -0.1282 & -0.0745 \\ 0.0846 & 0.9163 & -0.3938 \end{bmatrix},$$

$$U_{\text{PARAFAC}} = \begin{bmatrix} 0.0841 & 0.3795 & 0.9162 \\ 0.9929 & -0.1282 & -0.0745 \\ 0.0846 & 0.9163 & -0.3938 \end{bmatrix}, \quad U_{\text{Tucker}} = \begin{bmatrix} 0.0627 & 0.3707 & 0.9266 \\ 0.9952 & -0.0937 & -0.0298 \\ 0.0758 & 0.9240 & -0.3748 \end{bmatrix}.$$

Observations are as follows:

1. The $U_{\text{inc}}$ and $U_{\text{PARAFAC}}$ are not orthogonal.

2. Compared with Tucker, LROAT gives better approximations to the vectors $u_i$'s:

$$\|U - U_{\text{LROAT}}\| = 0.0252, \quad \|U - U_{\text{Tucker}}\| = 0.0527.$$

3. In terms of approximation quality, the residual norms (in percentage of the norm of $\mathcal{A}_5$), are

$$res_{\text{LROAT}} = 3.07\%, \quad res_{\text{inc}} = 1.36\%, \quad res_{\text{PARAFAC}} = 1.36\%, \quad res_{\text{Tucker}} = 0\%.$$

**7. Concluding remarks.** In the present paper we studied the tensor SVD, and characterized its existence in relation to the HOSVD. Similar to the concept of rank, the SVD of higher order tensors exhibits a quite different behavior and characteristics from those of matrices. Thus, the SVD of a matrix is guaranteed to exist, though it may have different representations due to orthogonal transformations of singular vectors corresponding to the same singular value. On the other hand, there are many ways in which a tensor can fail to have an SVD (see the results in Section 4), but when it exists, this decomposition is unique up to signs.

We have also discussed a new form of optimal low rank approximation of tensors, where orthogonality is required. This approximation is inspired by the constraints of the Tucker model and the PARAFAC model. In some applications, the proposed approximation model may be favored, since it results in $N$ sets of orthonormal vectors or, equivalently, $r$ mutually orthogonal unit rank-1 tensors with different weights. Among the advantages of this approximation over the Tucker model is the fact that it requires far fewer entries to represent the core, and that it is easier to interpret. Also, compared with the PARAFAC model, the orthogonality of vectors may be useful in some cases. Further, the LROAT algorithm for computing the proposed approximation does not seem to exhibit the well-known slow convergence from which the ALS algorithm for PARAFAC suffers.

A major restriction of the proposed model is that the number of terms $r$ can not exceed the smallest dimension of all modes of the tensor. A consequence is that the approximation may still be very different from the original tensor even when the maximum $r$ is employed. However we note that when performing data analysis, the interpretation of the vectors and the core tensor might be more important than how much is lost when the data is approximated.

A nice aspect of the proposed approximation is that the optimum of the objective function can theoretically be attained, in contrast to the PARAFAC model which is ill-posed in a strict mathematical sense. We presented an algorithm to compute this approximation, but the computed result is only optimal in a local neighborhood. It will be interesting to study for what tensors or what initial guesses the LROAT algorithm converges to the global optimum, or to devise a new algorithm to solve this optimization problem. It is an open problem how fast LROAT converges, although empirically convergence is observed to be linear. We also discussed the symmetric variant of LROAT and pointed out the possibility of its non-convergence. Hence the convergence properties of this variant, and the observed phenomenon that the original LROAT algorithm can yield same side-matrices for symmetric tensors, remain to be investigated.

**Appendix. Does the ALS algorithm for PARAFAC converge?** It has been pointed out that the ALS algorithm for computing the PARAFAC model may

converge very slowly due to degenerate solutions or multicollinearities, and many alternatives have been proposed to address this problem [36, 37, 26]. During iterations, the objective value monotonically decreases by the nature of the alternating least squares procedure, and since the sequence is bounded, it converges. However, a proof of the convergence of the parallel factors is lacking. In general it is assumed that these factors converge, but may take a very large number of iterations. In this section, we discuss an experiment showing that the general concept of convergence is unclear in this context. Though only one example is given, we note that the exhibited behavior is not rare for randomly generated tensors. (On the other hand it may be argued that tensors in real applications are far from being filled with random entries.)

We generate an order-3 tensor $\mathcal{A} \in \mathbb{R}^{3 \times 3 \times 3}$ and run the ALS algorithm on $r = 2$, i.e., to compute the approximation

$$\mathcal{A} \approx \lambda_1 u_1^{(1)} \otimes u_1^{(2)} \otimes u_1^{(3)} + \lambda_2 u_2^{(1)} \otimes u_2^{(2)} \otimes u_2^{(3)}.$$

The Matlab code which generates the tensor $\mathcal{A}$ is as follows:

```
A(:,:,1) = [.99 .29 .08; .44 .69 .19; .00 .49 .97];
A(:,:,2) = [.36 .64 .10; .13 .73 .89; .01 .02 .76];
A(:,:,3) = [.58 .55 .98; .68 .77 .04; .96 .61 .98];
```

We use $u_i^{(n)} = e_i$, where $n = 1, 2, 3$ and $e_i$ is the $i$-th column of the identity matrix, as the initial guess.

Denote $U^{(n)} = \left[ u_1^{(n)}, u_2^{(n)} \right]$ for $n = 1, 2, 3$. Two plots are shown after running $10^5$ iterations (see Figure A.1). Figure A.1(a) shows the "convergence" history for each $U^{(n)}$. The curves represent $\left\| U_{(p)}^{(n)} - U_{(p-1)}^{(n)} \right\|$, where $p$ is the index of the iterations. A necessary condition for convergence to occur is that all the three curves decrease to zero. However we see from the figure that this may not be the case. To test the conjecture that each of the curves tends to a nonzero value, we use the following expression

$$\log_{10} y = \frac{a}{(10^{-4} x)^{1/\alpha}} + b$$

to fit the tailing part of the curves (starting from the $2 \times 10^4$-th iteration). Table A.1 gives the fitting results for different $\alpha$'s. When the number of iterations tends to infinity, the value $10^b$ will show the limit of the differences between two consecutive $U^{(n)}$'s.

It is still difficult to conclude for this example that the iterations do not converge since rounding has not been taken into account. However, it makes no practical difference for this case whether the sequence actually converges or whether it is exceedingly slow to converge. The result, if convergence holds, will be an inordinate number of iterations to reach a desirable level of convergence, and the cost will be too high in practice. This can be made evident by examining Figure A.1(b), which plots the parallel factor $u_2^{(1)}$ over all iterations: The 3rd entry of $u_2^{(1)}$ decreases from 0.2540 at the $5 \times 10^4$-th iteration to 0.2517 at the $10^5$-th iteration.

REFERENCES

Table A.1

*Curve fitting for the three curves in Figure A.1(a) using different $\alpha$ values. The error is measured as the quadratic mean of fitting errors in logarithmic scales, i.e., the RMS of $|\log_{10} y - \log_{10} y_{fit}|$. It will be easier to understand this error by noticing that the vertical axis of Figure A.1(a) has a length 8 (after taking logarithm).*

| | $\alpha$ | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| | $a$ | 2.8407 | 2.7925 | 3.2657 | 3.8359 | 4.4404 |
| $U^{(1)}$ | $b$ | $-7.5060$ | $-8.1549$ | $-8.8047$ | $-9.4544$ | $-10.1042$ |
| | error | 0.0595 | 0.0301 | 0.0201 | 0.0151 | 0.0121 |
| | $a$ | 2.8408 | 2.7926 | 3.2658 | 3.8360 | 4.4406 |
| $U^{(2)}$ | $b$ | $-7.6428$ | $-8.2917$ | $-8.9415$ | $-9.5913$ | $-10.2411$ |
| | error | 0.0595 | 0.0301 | 0.0201 | 0.0151 | 0.0121 |
| | $a$ | 2.8411 | 2.7930 | 3.2662 | 3.8365 | 4.4411 |
| $U^{(3)}$ | $b$ | $-7.7219$ | $-8.3709$ | $-9.0208$ | $-9.6707$ | $-10.3205$ |
| | error | 0.0595 | 0.0301 | 0.0201 | 0.0151 | 0.0121 |

[1] E. ACAR, S. A. CAMTEPE, M. KRISHNAMOORTHY, AND B. YENER, *Modeling and multiway analysis of chatroom tensors*, in Proc. of IEEE Int. Conf. on Intelligence and Security Informatics (ISI 05), 2005.

[2] B. W. BADER AND T. G. KOLDA, *Algorithm 862: MATLAB tensor classes for fast algorithm prototyping*, ACM Trans. Math. Softw., 32 (2006).

[3] J. BERGE, J. DE LEEUW, AND P. M. KROONENBERG, *Some additional results on principal components analysis of three-mode data by means of alternating least squares algorithms*, Psychometrika, 52 (1987), pp. 183–191.

[4] R. BRO, *Review on multiway analysis in chemistry—2000-2005*, Critical Reviews in Analytical Chemistry, 36 (2006), pp. 279–293.

[5] J.-F. CARDOSO AND P. COMON, *Independent component analysis, a survey of some algebraic methods*, in Proc. of IEEE International Symposium on Circuits and Systems (ISCAS 96), 1996.

[6] J. D. CARROLL AND J.-J. CHANG, *Analysis of individual differences in multidimensional scaling via an n-way generalization of "Eckart-Young" decomposition*, Psychometrika, 35 (1970), pp. 283–319.

[7] P. COMON, *Independent component analysis, a new concept?*, Signal Processing, 36 (1994), pp. 287–314.

[8] ———, *Tensor decompositions*, in Mathematics of Signal Processing V, J. G. McWhirter and I. K. Proudler, eds., Oxford University Press, 2002.

[9] P. COMON, G. GOLUB, L.-H. LIM, AND B. MOURRAIN, *Symmetric tensor and symmetric tensor rank*, SIAM J. Matrix Anal. Appl., (to appear).

[10] L. DE LATHAUWER, *Signal Processing based on Multilinear Algebra*, PhD thesis, Katholieke Universiteit Leuven, 1997.

[11] ———, *A link between the canonical decomposition in multilinear algebra and simultaneous matrix diagonalization*, SIAM J. Matrix Anal. Appl., 28 (2006), pp. 642–666.

[12] L. DE LATHAUWER, B. DE MOOR, AND J. VANDEWALLE, *A multilinear singular value decomposition*, SIAM J. Matrix Anal. Appl., 21 (2000), pp. 1253–1278.

[13] ———, *On the best rank-1 and rank-$(R_1, R_2, \ldots, R_N)$ approximation of higher-order tensors*, SIAM J. Matrix Anal. Appl., 21 (2000), pp. 1324–1342.

[14] ———, *Independent component analysis and (simultaneous) third-order tensor diagonalization*, IEEE Trans. Signal Process., 49 (2001), pp. 2262–2271.

[15] ———, *Computation of the canonical decomposition by means of a simultaneous generalized schur decomposition*, SIAM J. Matrix Anal. Appl., 26 (2004), pp. 295–327.

[16] B. DE MOOR, *Daisy: Database for the identification of systems.* URL: `http://homes.esat.kuleuven.be/~smc/daisy/`. Used dataset: `tongue.dat`, section: Biomedical Systems, code: 97-001.

[17] V. DE SILVA AND L.-H. LIM, *Tensor rank and the ill-posedness of the best low-rank approximation problem*, SIAM J. Matrix Anal. Appl., (to appear).

[18] C. ECKART AND G. YOUNG, *The approximation of one matrix by another of lower rank*, Psychometrika, 1 (1936), pp. 211–218.

[19] R. A. Harshman, *Foundations of the PARAFAC procedure: Models and conditions for an "explanatory" multimodal factor analysis*, UCLA Working Papers in Phonetics, 16 (1970), pp. 1–84.

[20] R. A. Harshman, P. Ladefoged, and L. Goldstein, *Factor analysis of tongue shapes*, J. Acoust. Soc. Am., 62 (1977), pp. 693–707.

[21] N. J. Higham and R. S. Schreiber, *Fast polar decomposition of an arbitrary matrix*, SIAM J. Sci. Comput., 11 (1990), pp. 648–655.

[22] R. A. Horn and C. R. Johnson, *Matrix Analysis*, Cambridge University Press, 1985.

[23] J. Jájá and J. Takche, *On the validity of the direct sum conjecture*, SIAM J. Comput., 15 (1986), pp. 1004–1020.

[24] H. A. L. Kiers, *TUCKALS core rotations and constrained TUCKALS modelling*, Statistica Applicata, 4 (1992), pp. 659–667.

[25] ———, *An alternating least squares algorithms for PARAFAC2 and three-way DEDICOM*, Comput Statist. Data Anal., 16 (1993), pp. 103–118.

[26] ———, *A three-step algorithm for CANDECOMP/PARAFAC analysis of large data sets with multicollinearity*, J. Chemometrics, 12 (1998), pp. 155–171.

[27] E. Kofidis and P. A. Regalia, *On the best rank-1 approximation of higher-order supersymmetric tensors*, SIAM J. Matrix Anal. Appl., 23 (2001), pp. 863–884.

[28] T. G. Kolda, *Orthogonal tensor decompositions*, SIAM J. Matrix Anal. Appl., 23 (2001), pp. 243–255.

[29] ———, *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 (2003), pp. 762–767.

[30] T. G. Kolda and B. W. Bader, *Tensor decompositions and applications*, SIAM Rev., (to appear).

[31] T. G. Kolda, B. W. Bader, and J. P. Kenny, *Higher-order web link analysis using multilinear algebra*, in Proc. of the 5th IEEE International Conference on Data Mining (ICDM 05), 2005.

[32] W. P. Krijnen, *Convergence of the sequence of parameters generated by alternating least squares algorithms*, Comput Statist. Data Anal., 51 (2006), pp. 481–489.

[33] P. M. Kroonenberg and J. De Leeuw, *Principal component analysis of three-mode data by means of alternating least squares algorithms*, Psychometrika, 45 (1980), pp. 69–97.

[34] S. E. Leurgans, R. T. Ross, and R. B. Abel, *A decomposition for three-way arrays*, SIAM J. Matrix Anal. Appl., 14 (1993), pp. 1064–1083.

[35] C. D. Martin and C. Van Loan, *A Jacobi-type method for computing orthogonal tensor decompositions*, SIAM J. Matrix Anal. Appl., (to appear).

[36] B. C. Mitchell and D. S. Burdick, *Slowly converging parafac sequences: Swamps and two-factor degeneracies*, J. Chemometrics, 8 (1994), pp. 155–168.

[37] P. Paatero, *A weighted non-negative least squares algorithm for three-way 'PARAFAC' factor analysis*, Chemometrics Intell. Lab. Syst., 38 (1997), pp. 223–242.

[38] A. Smilde, R. Bro, and P. Geladi, *Multi-way Analysis: Applications in the Chemical Sciences*, Wiley, 2004.

[39] J. Sun and C. Chen, *Generalized polar decomposition*, Math. Numer. Sin., 11 (1989), pp. 262–273.

[40] L. R. Tucker, *Some mathematical notes on three-mode factor analysis*, Psychometrika, 31 (1966), pp. 279–311.

[41] M. A. O. Vasilescu and D. Terzopoulos, *Multilinear subspace analysis for image ensembles*, in Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR 03), 2003.

[42] H. Wang, Q. Wu, L. Shi, Y. Yu, and N. Ahuja, *Out-of-core tensor approximation of multidimensional matrices of visual data*, ACM Trans. Gr., 24 (2005), pp. 527–535.

[43] T. Zhang and G. H. Golub, *Rank-one approximation to high order tensors*, SIAM J. Matrix Anal. Appl., 23 (2001), pp. 534–550.
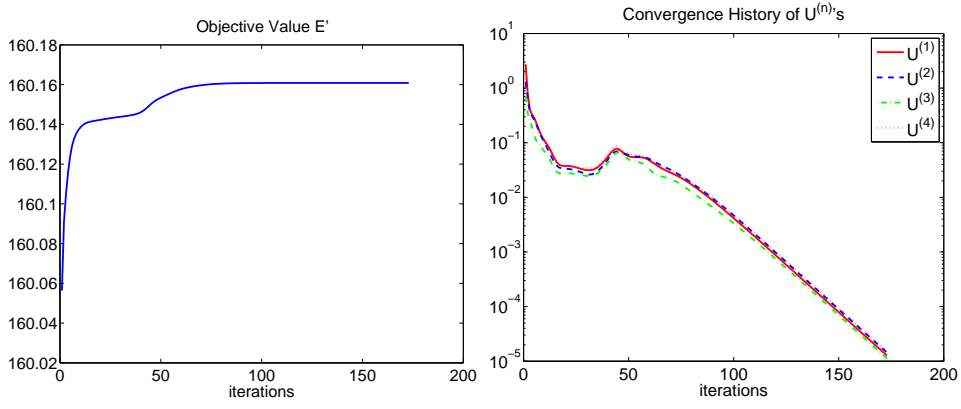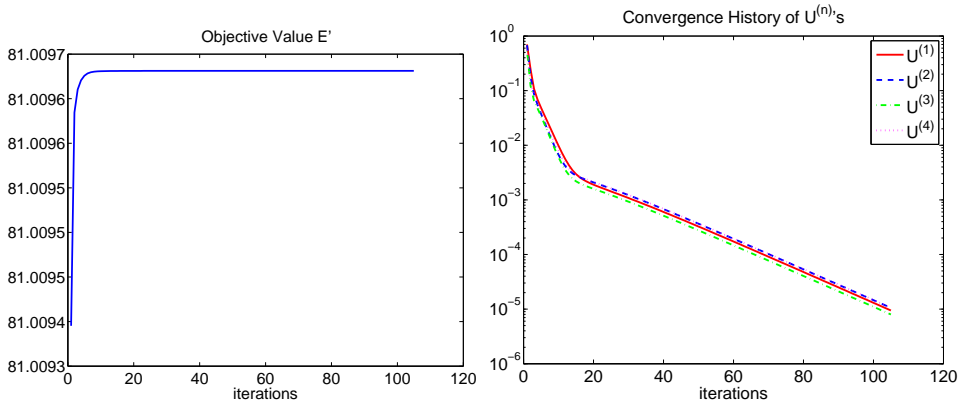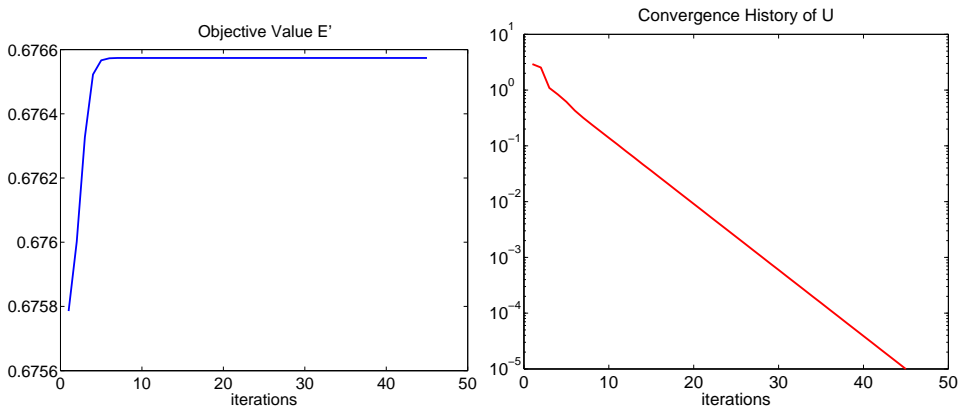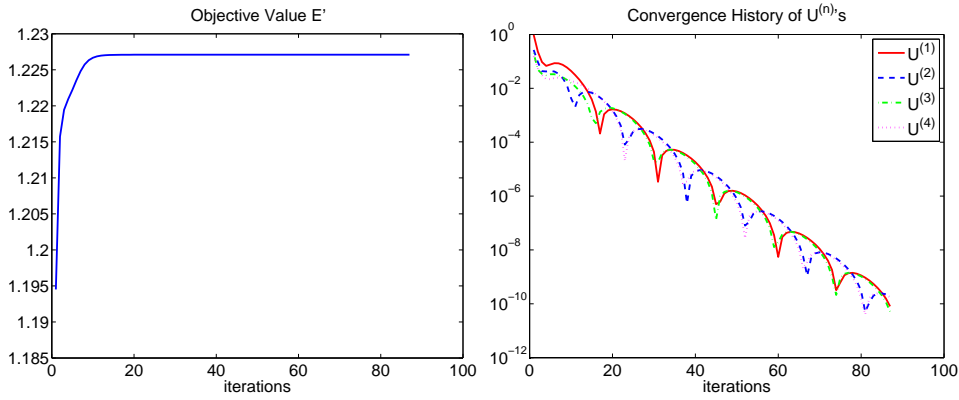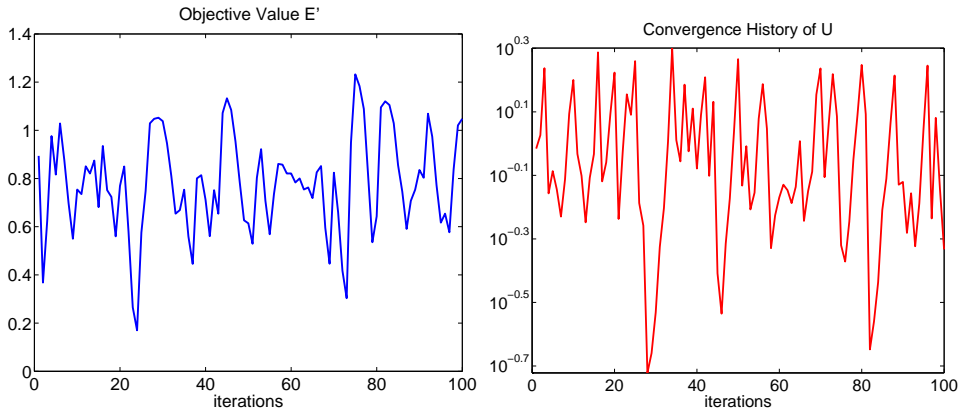
(a) Tensor $\mathcal{A}_1$: Randomly generated.



(b) Tensor $\mathcal{A}_2$: Low rank plus Gaussian noise.



(c) Tensor $\mathcal{A}_3$: $(\mathcal{A}_3)_{ijk} = 1/(i^2 + j^2 + k^2)$.
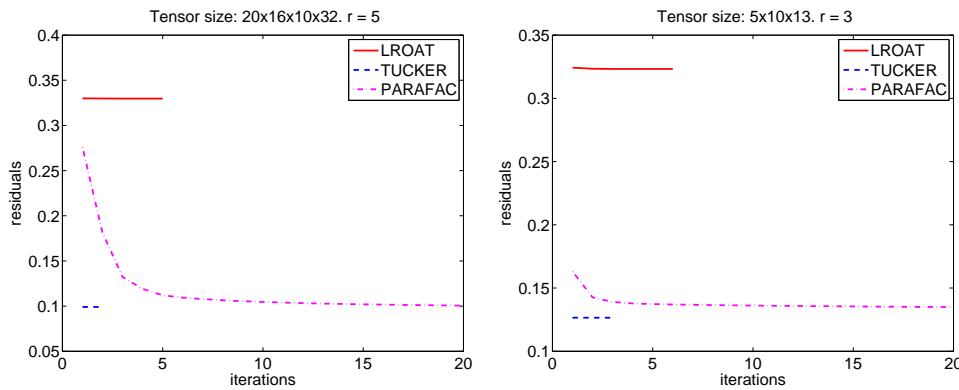
FIG. 6.1. *Experiment 1: Convergence tests for LROAT.*

(a) Tensor $\mathcal{A}_4$ directly applied to LROAT.



(b) Tensor $\mathcal{A}_4$ applied to the symmetric variant of LROAT.
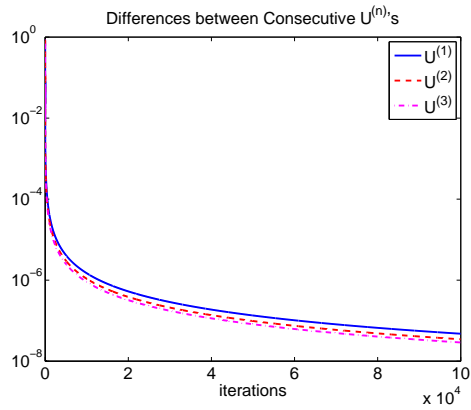
FIG. 6.2. *Experiment 1 (continued): Convergence tests for LROAT.*
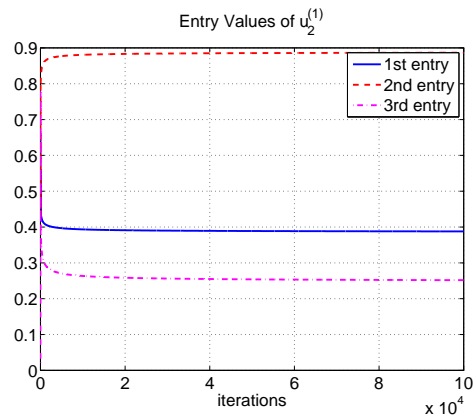


(a) Low-rank-plus-Gaussian-noise tensor.                    (b) Real-life tensor.

FIG. 6.3. *Experiment 2: Comparison of LROAT, Tucker and PARAFAC.*

(a) Differences of the $U^{(n)}$'s between consecutive iterations.

(b) Entry values of the parallel factor $u_2^{(1)}$.

FIG. A.1. *Slow convergence or non-convergence of ALS for PARAFAC.*