

# Efficient Alternating Least Squares Algorithms for Low Multilinear Rank Approximation of Tensors

Chuanfu Xiao · Chao Yang · Min Li

Received: date / Accepted: date

**Abstract** The low multilinear rank approximation, also known as the truncated Tucker decomposition, has been extensively utilized in many applications that involve higher-order tensors. Popular methods for low multilinear rank approximation usually rely directly on matrix SVD, therefore often suffer from the notorious intermediate data explosion issue and are not easy to parallelize, especially when the input tensor is large. In this paper, we propose a new class of truncated HOSVD algorithms based on alternating least squares (ALS) for efficiently computing the low multilinear rank approximation of tensors. The proposed ALS-based approaches are able to eliminate the redundant computations of the singular vectors of intermediate matrices and are therefore free of data explosion. Also, the new methods are more flexible with adjustable convergence tolerance and are intrinsically parallelizable on high-performance computers. Theoretical analysis reveals that the ALS iteration in the proposed algorithms is  $q$ -linear convergent with a relatively wide convergence region. Numerical experiments with large-scale tensors from both synthetic and real-world applications demonstrate that ALS-based methods can substantially reduce the total cost of the original ones and are highly scalable for parallel computing.

**Keywords** Low multilinear rank approximation · Truncated Tucker decomposition · Alternating least squares · Parallelization

**Mathematics Subject Classification (2010)** 15A69 · 49M27 · 65D15 · 65F55

---

Chuanfu Xiao  
School of Mathematical Sciences, Peking University, Beijing 100871, China  
E-mail: chuanfuxiao@pku.edu.cn

Chao Yang (✉)  
School of Mathematical Sciences, Peking University, Beijing 100871, China  
E-mail: chao.yang@pku.edu.cn

Min Li  
Institute of Software, Chinese Academy of Sciences, Beijing 100190, China  
E-mail: limin2016@iscas.ac.cn

## 1 Introduction

As a natural extension of vectors (first-order) and matrices (second-order), higher-order tensors have been receiving increasingly more attention in various applications, such as signal processing [17, 56], computer vision [55, 64, 66], chemometrics [27, 33, 9], deep learning [13, 45], and scientific computing [8, 34, 24]. For decades, tensor decompositions have been extensively utilized as an efficient tool for dimension reductions, latent variable analysis and other purposes in a wide range of scientific and engineering fields [24, 7, 37, 31, 29, 20]. There exist a number of tensor decomposition models, such as canonical polyadic (CP, or CANDECOMP/PARAFAC) decomposition [28, 12, 35], Tucker decomposition [61, 41, 15], tensor train (TT) model [49], and hierarchical Tucker (HT) model [23, 25, 48]. Among them, the Tucker decomposition, also known as the higher-order singular value decomposition (HOSVD), is regarded as a generalization of the matrix singular value decomposition (SVD) and has been applied with significant successes in many applications [61, 41, 15, 17, 31, 36].

In both theory and practice, a commonly considered tensor computation problem is the low multilinear rank approximation [16, 36], also known as the truncated Tucker decomposition [63, 4], which reads

$$\min_{\mathcal{B}} \|\mathcal{A} - \mathcal{B}\|, \quad (1.1)$$

where  $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$  is a given  $N$ th-order tensor and  $\mathcal{B} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$  is its low multilinear rank approximation [15, 36]. Existing approaches for solving (1.1) can be roughly divided into two categories [62]: *non-iterative* and *iterative* methods. The most popular *non-iterative* algorithms for the low multilinear rank approximation of higher-order tensors is the truncated HOSVD ( $t$ -HOSVD) [61, 16] and its improved version, the sequentially truncated HOSVD ( $st$ -HOSVD) [63]. Despite the fact that the results of  $t$ -HOSVD and  $st$ -HOSVD are usually suboptimal, they can serve as good initial solution for popular iterative methods such as higher-order orthogonal iteration (HOOI) [38, 16]. Other than the HOOI method, which is a first-order iterative method, some efforts are also made in developing second-order approaches, such as Newton-type [21, 53, 32] and trust-region [31, 30] algorithms. Although these methods can achieve faster convergence under certain conditions, they are still in early study and are usually not suitable for large-scale tensors [70].

In this paper, we focus on studying how to efficiently compute the truncated Tucker decomposition (1.1) of higher-order tensors by modifying the  $t$ -HOSVD and  $st$ -HOSVD algorithms so that certain fast iterative procedures are included. As a major cost of the two algorithms, the computation of tensor-matrix multiplications has been extensively optimized in a number of high-performance tensor libraries [50, 42, 57, 54], and therefore is not the focus of this study. Another potential bottleneck of the  $t$ -HOSVD and  $st$ -HOSVD algorithms is the calculation of the singular vectors of the intermediate matrices, which can be done by applying SVD to the matricized tensor or eigen-decomposition to the Gram matrix [38, 16, 36, 4, 46]. The matrix SVD can be obtained by using Krylov subspace methods [22, 14, 6], whilst the eigen-decomposition of the symmetric nonnegative definite Gram matrix can be done with a Krylov-Schur algorithm [58, 59, 69, 22]. Due to the fact that these methods rely on the factorization of intermediate matrices, they suffer from the notorious *data explosion* issue [4, 46, 47]. And even if the hardware storage allowed, they are still not scalable for parallel computing and the total computation cost could be unbearably high.

In order to improve the performance of the  $t$ -HOSVD and  $st$ -HOSVD algorithms, we propose a class of alternating least squares (ALS) based algorithms for efficiently calculating the low multilinear rank approximation of tensors. The key observation is that in the original algorithms the computations of singular vectors of the intermediate matrices are indeed not necessary and can be

replaced with low rank approximations, and the low rank approximations can be done by using an ALS method which does not explicitly require intermediate tensor matricization, with the help of a row-wise update rule. The proposed ALS-based algorithms enjoy advantages in computing efficiency, error adaptivity and parallel scalability, especially for large-scale tensors. We present theoretical analysis and show that the ALS iteration in the proposed algorithms is q-linear convergent with a relatively wide convergence region. Several numerical experiments with both synthetic and real-world tensor data demonstrate that new algorithms can effectively alleviate the data explosion issue of the original ones and are highly parallelizable on parallel computers.

The organization of the paper is as follows. In Sec. 2, we introduce some basic notations of tensor and the corresponding algorithms. In Sec. 3, the  $t$ -HOSVD-ALS and  $st$ -HOSVD-ALS algorithms are proposed. Some theoretical analysis on the convergence behavior of the ALS methods can also be found in Sec. 3. After that, computational complexity and the approximation errors of proposed algorithms are analyzed in Sec. 4. Test results on several numerical experiments are reported in Sec. 5. And the paper is concluded in Sec. 6.

## 2 Notations and Nomenclatures

Symbols frequently used in this paper can be found in the following table.

Symbols	Notations
$a$	scalar
$\mathbf{a}$	vector
$\mathbf{A}$	matrix
$\mathcal{A}$	three or higher-order tensor
$\circ$	vector outer product
$\times_n$	mode- $n$ product of tensor and matrix
$\mathbf{I}_n$	identity matrix with size $n \times n$
$I_{n:N}$	$\prod_{i=n}^N I_i$
$\mathcal{R}(\mathbf{A})$	a subspace formed by the columns of matrix $\mathbf{A}$
$\sigma(\mathbf{A})$	a set that consists of singular values of matrix $\mathbf{A}$
$\mathbf{A}^\dagger$	pseudo-inverse of matrix $\mathbf{A}$

Given an  $N$ th-order tensor  $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_N}$ , we denote  $\mathcal{A}_{i_1, i_2, \dots, i_N}$  as its  $(i_1, i_2, \dots, i_N)$ -th element. In particular, rank one tensor is denoted as

$$\mathbf{u}_1 \circ \mathbf{u}_2 \circ \cdots \circ \mathbf{u}_N,$$

where  $\mathbf{u}_n \in \mathbb{R}^{I_n}$  is a vector.

The Frobenius norm of tensor  $\mathcal{A}$  is defined as

$$\|\mathcal{A}\|_F = \sqrt{\sum_{i_1, i_2, \dots, i_N} \mathcal{A}_{i_1, i_2, \dots, i_N}^2}.$$

The matricization of a higher-order tensor is a process of reordering the elements of the tensor into a matrix. For example, the mode- $n$  matricization of tensor  $\mathcal{A}$  is denoted as  $\mathbf{A}_{(n)}$ , which is a matrix belonging to  $\mathbb{R}^{I_n \times (I_1 \cdots I_{n-1} I_{n+1} \cdots I_N)}$ . Specifically, the  $(i_1, i_2, \dots, i_N)$ -th element of tensor  $\mathcal{A}$  is mapped to the  $(i_n, j)$ -th entry of matrix  $\mathbf{A}_{(n)}$ , where

$$j = 1 + \sum_{k=1, k \neq n}^N (i_k - 1) J_k \quad \text{with} \quad J_k = \prod_{m=1, m \neq n}^{k-1} I_m.$$

The multilinear rank of a higher-order tensor  $\mathcal{A}$  is an integer array  $(R_1, R_2, \dots, R_N)$ , where  $R_n$  is the rank of its mode- $n$  matricization  $\mathbf{A}_{(n)}$ .

A frequently encountered operation in tensor computation is the tensor-matrix multiplication. In particular, the mode- $n$  tensor-matrix multiplication refers to the contraction of the tensor with a matrix along the  $n$ -th index. For example, suppose that  $\mathbf{U} \in \mathbb{R}^{J \times I_n}$  is a matrix, the mode- $n$  product of  $\mathcal{A}$  and  $\mathbf{U}$  is denoted as  $\mathcal{A} \times_n \mathbf{U} \in \mathbb{R}^{I_1 \times \cdots \times I_{n-1} \times J \times I_{n+1} \times \cdots \times I_N}$ . Elementwisely, one has

$$\mathcal{B}_{i_1, \dots, j, \dots, i_N} = (\mathcal{A} \times_n \mathbf{U})_{i_1, \dots, j, \dots, i_N} = \sum_{i_n=1}^{I_n} \mathcal{A}_{i_1, \dots, i_n, \dots, i_N} \mathbf{U}_{j, i_n}.$$

The Tucker decomposition [61, 41], also known as the higher-order singular value decomposition (HOSVD) [15], is formally defined as

$$\mathcal{A} = \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \cdots \times_N \mathbf{U}^{(N)},$$

where  $\mathcal{G} \in \mathbb{R}^{R_1 \times R_2 \times \cdots \times R_N}$  is referred to as the core tensor and  $\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times R_n}$  are column orthogonal with each other for all  $n \in \{1, 2, \dots, N\}$ . We remark here that the size of the core tensor is often smaller than that of the original tensor, though it is hard to know how small it can be a priori [19, 36]. In many applications, the Tucker decomposition is usually applied in its truncated form, which reads

$$\begin{aligned} \min_{\mathcal{G}; \mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \dots, \mathbf{U}^{(N)}} \|\mathcal{A} - \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \cdots \times_N \mathbf{U}^{(N)}\|, \\ \text{s.t.} \quad \mathbf{U}^{(n)T} \mathbf{U}^{(n)} = \mathbf{I}_{R_n}, \quad n \in \{1, 2, \dots, N\} \end{aligned} \quad (2.1)$$

where  $(R_1, R_2, \dots, R_N)$  is a pre-determined truncation, smaller than the size of original tensor. Suppose that the exact solution of (2.1) is  $\mathbf{U}^{*(1)}, \mathbf{U}^{*(2)}, \dots, \mathbf{U}^{*(N)}$ , and  $\mathcal{G}^*$ , then it is easy to see that

$$\mathcal{G}^* = \mathcal{A} \times_1 (\mathbf{U}^{*(1)})^T \times_2 (\mathbf{U}^{*(2)})^T \cdots \times_N (\mathbf{U}^{*(N)})^T,$$

which means

$$\mathcal{A} \times_1 (\mathbf{U}^{*(1)}) (\mathbf{U}^{*(1)})^T \times_2 (\mathbf{U}^{*(2)}) (\mathbf{U}^{*(2)})^T \cdots \times_N (\mathbf{U}^{*(N)}) (\mathbf{U}^{*(N)})^T \quad (2.2)$$

is the best low multilinear rank approximation of  $\mathcal{A}$ .

To compute the best low multilinear rank approximation of a higher-order tensor in the truncated Tucker decomposition, a popular approach is the truncated HOSVD ( $t$ -HOSVD, [61]) originally presented by Tucker himself [61]. Nowadays, it is better known with the effort of Lathauwer *et al.* [16], who analyzed the structure of core tensor and proposed to employ SVD of the intermediate matrices in truncated HOSVD. The computing procedure of  $t$ -HOSVD is given in Algorithm 1.

**Algorithm 1**  $t$ -HOSVD [61,16]

---

**Input:** Tensor  $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ , truncation  $(R_1, R_2, \dots, R_N)$   
**Output:** Low multilinear rank approximation  $\hat{\mathcal{A}} \approx \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \dots \times_N \mathbf{U}^{(N)}$   
1: **for all**  $n \in \{1, 2, \dots, N\}$  **do**  
2:    $\mathbf{Q} \leftarrow$  leading left singular vectors of  $\mathbf{A}_{(n)}$   
3:    $\mathbf{U}^{(n)} \leftarrow \mathbf{Q}$   
4: **end for**  
5:  $\mathcal{G} \leftarrow \mathcal{A} \times_1 \mathbf{U}^{(1)T} \times_2 \mathbf{U}^{(2)T} \dots \times_N \mathbf{U}^{(N)T}$

---

We remark here that in Algorithm 1, the computation of  $\mathbf{Q}$  can also be done by calculating the  $R_n$  eigenvectors of the Gram matrix  $\mathbf{A}_{(n)}\mathbf{A}_{(n)}^T$ . It is clear that  $t$ -HOSVD can be seen as a natural extension of the truncated SVD of a matrix to higher-order tensors. But unlike the matrix case, the approximation error of  $t$ -HOSVD is quasi-optimal [61,16,36,63].

As a subsequent improvement of  $t$ -HOSVD, the sequentially truncated HOSVD ( $st$ -HOSVD), proposed by Vannieuwenhoven *et al.* [63] uses a different truncation strategy, as shown in Algorithm 2.

**Algorithm 2**  $st$ -HOSVD [63]

---

**Input:** Tensor  $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ , truncation  $(R_1, R_2, \dots, R_N)$   
**Output:** Low multilinear rank approximation  $\hat{\mathcal{A}} \approx \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \dots \times_N \mathbf{U}^{(N)}$   
1: Select an order of  $\{1, 2, \dots, N\}$ , i.e.,  $\{i_1, i_2, \dots, i_N\}$ .  
2:  $\mathcal{B} \leftarrow \mathcal{A}$   
3: **for all**  $n \in \{i_1, i_2, \dots, i_N\}$  **do**  
4:    $\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}^T \leftarrow$  matrix SVD of  $\mathbf{B}_{(n)}$   
5:    $\mathbf{U}^{(n)} \leftarrow \mathbf{U}$   
6:    $\mathcal{B} \leftarrow \mathbf{\Sigma}\mathbf{V}^T$  in tensor format  
7: **end for**  
8:  $\mathcal{G} \leftarrow \mathcal{B}$

---

Analogous to  $t$ -HOSVD, in Algorithm 2 one can also obtain  $\mathbf{Q}$  by computing the  $R_n$  eigenvectors of the Gram matrix  $\mathbf{B}_{(n)}\mathbf{B}_{(n)}^T$ , and the core tensor is updated by  $\mathcal{B} = \mathcal{B} \times_n \mathbf{U}^{(n)T}$ .

Unlike  $t$ -HOSVD, the  $st$ -HOSVD algorithm does not compute the singular vectors of  $\mathbf{A}_{(n)}$ . In particular, the factor matrices and core tensor in  $st$ -HOSVD are calculated simultaneously, which greatly reduces the size of the intermediate matrices. It is worth mentioning that the order of  $\{1, 2, \dots, N\}$  in Algorithm 2 could have a strong influence on the computational cost and approximation error of  $st$ -HOSVD. But there is no theoretical guidance on how to select the order. A heuristic suggestion on how to select the processing order of  $st$ -HOSVD was given based on the dimension of each mode [63], i.e.,  $I_n$ ,  $n = 1, 2, \dots, N$ . As a supplement, the order can also be decided according to the truncation  $(R_1, R_2, \dots, R_N)$  of the tensor, which is summarized in the following proposition.

**Proposition 1** *Let  $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$  be an  $N$ th-order tensor, and the truncation is set to  $(R_1, R_2, \dots, R_N)$ . Without loss of generality, suppose that  $I_n \approx I$  for any  $n \in \{1, 2, \dots, N\}$ , and*

$R_1 \leq R_2 \leq \dots \leq R_N$ . Then the *st*-HOSVD algorithm based on order  $\{1, 2, \dots, N\}$  has the lowest computational cost, as compared with other computational orders.

**Proof.** If we select  $\{1, 2, \dots, N\}$  as the order of Algorithm 2, when applying a Krylov subspace method to compute the truncated matrix SVD, the computational cost is

$$\mathcal{O}\left(\sum_{n=1}^N R_{1:n} I_{n:N}\right) \approx \mathcal{O}\left(\sum_{n=1}^N R_{1:n} I^{N-n+1}\right), \quad (2.3)$$

Similarly, the computational cost when we select  $(i_1, i_2, \dots, i_N)$  as the order of Algorithm 2 is

$$\mathcal{O}\left(\sum_{n=1}^N R_{i_1:i_n} I_{i_n:i_N}\right) \approx \mathcal{O}\left(\sum_{n=1}^N R_{i_1:i_n} I^{N-n+1}\right). \quad (2.4)$$

Clearly, (2.3) is smaller than (2.4).  $\square$

For the best low multilinear rank approximation (2.1), it is easy to see is that  $\mathbf{U}^{*(n)}$  is a column orthogonal factor matrix, therefore  $(\mathbf{U}^{*(n)})(\mathbf{U}^{*(n)})^T$  represents the orthogonal projection of subspace  $\mathcal{R}(\mathbf{U}^{*(n)})$ . Consequently, subspace represented by the optimal factor matrices are critical. SVD and eigen-decomposition are the commonly applied approaches to determine this subspace in the original *t*-HOSVD and *st*-HOSVD procedures. However, because of the introduction of the intermediate matrices, both SVD and eigen-decomposition suffer from the notorious data explosion issue. Although some efforts have been made to alleviate the data explosion problem by, e.g., introducing an implicit Arnoldi procedure, these fixes are usually not generalizable to large-scale tensors in real applications and are not parallelization friendly.

In addition to SVD or eigen-decomposition, tensor matricization and tensor-matrix multiplication are also important in the original *t*-HOSVD and *st*-HOSVD algorithms. Recently, some efforts on high-performance optimizations of basic tensor operations are made. For example, Li *et al.* proposed a shared-memory parallel implementation of dense tensor-matrix multiplication [42], and Smith *et al.* considered sparse tensor-matrix multiplications [57]. Nevertheless, the calculation of SVD or eigen-decomposition is still the major challenge in the *t*-HOSVD and *st*-HOSVD algorithms, especially for large-scale tensors.

### 3 Alternating Least Squares Algorithms for *t*-HOSVD and *st*-HOSVD

In this paper we tackle the challenges of the original *t*-HOSVD and *st*-HOSVD algorithms from an alternating least squares (ALS) perspective. Instead of utilizing SVD or eigen-decomposition on the intermediate matrices, we propose to compute the dominant subspace with an ALS method to solve a closely related matrix low rank approximation problem. The classical ALS method for solving matrix low rank approximation problems was originally proposed by Leeuw et al., [18] and further applied in principal component analysis [71]. Algorithm 3 shows the detailed procedure of the ALS method.

**Algorithm 3**  $[\mathbf{L}^*, \mathbf{R}^*] = \text{ALS}(\mathbf{A}, r)$ **Input:** Matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , truncation  $r < \min\{m, n\}$ Initial guesses  $\mathbf{L}_0 \in \mathbb{R}^{m \times r}$  or  $\mathbf{R}_0 \in \mathbb{R}^{n \times r}$ **Output:** Low rank approximation  $\hat{\mathbf{A}} = \mathbf{L}^* \mathbf{R}^{*T}$ 

- 1:  $k \leftarrow 0$
- 2: **while** not convergent **do**
- 3: Solving multi-side least squares problem  $\min_{\mathbf{R}} \|\mathbf{L}_k \mathbf{R}^T - \mathbf{A}\|_F^2$
- 4:  $\mathbf{R}_k \leftarrow (\mathbf{A}^T \mathbf{L}_k)(\mathbf{L}_k^T \mathbf{L}_k)^{-1}$
- 5: Solving multi-side least squares problem  $\min_{\mathbf{L}} \|\mathbf{R}_k \mathbf{L}^T - \mathbf{A}^T\|_F^2$
- 6:  $\mathbf{L}_{k+1} \leftarrow (\mathbf{A} \mathbf{R}_k)(\mathbf{R}_k^T \mathbf{R}_k)^{-1}$
- 7:  $k \leftarrow k + 1$
- 8: **end while**

**Remark 1** We remark that unlike using an iterative method to solve matrix singular value problems, as was suggested in [63] to replace the matrix SVD, the proposed ALS-based approach can avoid the computations of singular pairs/triplets and thus achieve higher performance.

**Remark 2** To consider the existence and uniqueness of the solution of the ALS iterations, we note that in line 3 of Algorithm 3, if the coefficient matrices  $\mathbf{L}_k$  are nonsingular, the multi-side least squares problem  $\min_{\mathbf{R}} \|\mathbf{L}_k \mathbf{R}^T - \mathbf{A}\|$  is equivalent to linear equation  $\mathbf{L}_k^T \mathbf{L}_k \mathbf{R}^T = \mathbf{L}_k^T \mathbf{A}$ , which has a unique solution as shown in line 4 of Algorithm 3. The consistency of nonsingularity is further interpreted in Remark 3.

As an iterative method, the number of iterations for the ALS method has a dependency on the initial guess and the convergence criterion [60]. In what follows we will establish a rigorous convergence theory of the ALS method and derive an evaluation of the convergence region, which can help understand how the initial guess could affect the speed of convergence.

To establish the convergence theory of the ALS method, we first require the following lemma, which was proved in [68].

**Lemma 1** Let  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$  be symmetric positive definite matrices and satisfy

$$\mathbf{B} \leq \mathbf{A},$$

then the following inequalities hold

$$\|\mathbf{A}^{-1} \mathbf{B}\|_2 \leq 1 \text{ and } \|\mathbf{B} \mathbf{A}^{-1}\|_2 \leq 1,$$

where  $\mathbf{B} \leq \mathbf{A}$  represents  $\mathbf{A} - \mathbf{B}$  is symmetric semi-positive matrix.

The convergence theorem of Algorithm 3 is summarized in the theorem below.

**Theorem 1** Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$  be a matrix, and  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min\{m, n\}}$  be the singular values. Suppose that the following conditions hold:

1°  $\sigma(\mathbf{L}_k), \sigma(\mathbf{R}_k)$  are uniformly bounded.

2°  $\mathcal{R}(\mathbf{L}_0)$  is in a neighborhood of the exact solution.

Then Algorithm 3 is local  $q$ -linear convergent, and the convergence ratio is approximately  $\sigma_{r+1}^2 / \sigma_r^2$ , where  $\sigma_{r+1} < \sigma_r$ .

This theorem illustrates the convergence of the ALS method in a viewpoint of subspace, and the convergence ratio depends on the gap of  $\sigma_r$  and  $\sigma_{r+1}$ . The detailed proof can be found in Appendix A.

**Remark 3** *If condition 1° in Theorem 1 is not satisfied, then either  $\mathbf{L}_k$  or  $\mathbf{R}_k$  is close to singular. This implies that the truncation  $r$  is inappropriately chosen, i.e., greater than the numerical rank of  $\mathbf{A}$ .*

An evaluation of the convergence region of the ALS method can be found in the following theorem.

**Theorem 2** *Under the assumption of Theorem 1, provided that the initial guess  $\mathbf{L}_0$  satisfies*

$$\|\mathbf{L}_0^{(2)}(\mathbf{L}_0^{(1)})^{-1}\|_2 \leq \sqrt{\frac{\sigma_r^2 - (\sigma_r - \varepsilon)^2}{(\sigma_r - \varepsilon)^2 - \sigma_{min}^2}}, \quad (3.1)$$

*then the ALS method converges to the exact solution. Here*

$$\mathbf{U}^T \mathbf{L}_0 = \begin{pmatrix} \mathbf{U}_1^T \mathbf{L}_0 \\ \mathbf{U}_2^T \mathbf{L}_0 \end{pmatrix} = \begin{pmatrix} \mathbf{L}_0^{(1)} \\ \mathbf{L}_0^{(2)} \end{pmatrix},$$

$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  is the full SVD of  $\mathbf{A}$ ,  $\mathbf{U} = [\mathbf{U}_1, \mathbf{U}_2]$  is the block form of  $\mathbf{U}$ , and  $\varepsilon$  is an arbitrary positive number such that

$$\sigma_r - \varepsilon > \sigma_{r+1}.$$

The proof of Theorem 2 is provided in Appendix B. We remark that it can be seen from the theorem that, within the convergence region, a better initial guess guarantees faster convergence. It is also worth noting that (3.1) indicates that the convergence region depends on  $\varepsilon$ . A smaller  $\varepsilon$  means higher requirement for the initial guess, but less number of iterations.

With the help of Algorithm 3, we are able to solve the rank- $R_n$  approximation problem to obtain the dominant subspace of  $\mathbf{A}_{(n)}$  in  $t$ -HOSVD. Based on it, we derive the ALS accelerated versions of the  $t$ -HOSVD algorithm, namely  $t$ -HOSVD-ALS, presented in Algorithm 4.

---

#### Algorithm 4 $t$ -HOSVD-ALS

---

**Input:** Tensor  $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ , truncation  $(R_1, R_2, \dots, R_N)$

**Output:** Low multilinear rank approximation  $\hat{\mathcal{A}} \approx \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \dots \times_N \mathbf{U}^{(N)}$

- 1: **for all**  $n \in \{1, 2, \dots, N\}$  **do**
  - 2:    $\mathbf{A}_{(n)} \leftarrow \mathcal{A}$  in matrix format
  - 3:    $\mathbf{L}, \mathbf{R} \leftarrow \text{ALS}(\mathbf{A}_{(n)}, R_n)$
  - 4:    $\hat{\mathbf{Q}}, \hat{\mathbf{R}} \leftarrow$  reduced QR decomposition of  $\mathbf{L}$
  - 5:    $\mathbf{U}^{(n)} \leftarrow \hat{\mathbf{Q}}$
  - 6: **end for**
  - 7:  $\mathcal{G} \leftarrow \mathcal{A} \times_1 \mathbf{U}^{(1)T} \times_2 \mathbf{U}^{(2)T} \dots \times_N \mathbf{U}^{(N)T}$
- 

Similar to  $st$ -HOSVD, the proposed  $t$ -HOSVD-ALS algorithm does not explicitly compute the singular vectors of  $\mathbf{A}_{(n)}$  either. Instead, only an orthogonal basis is computed for the dominant subspace of  $\mathbf{A}_{(n)}$  in  $t$ -HOSVD-ALS. It is obtained by QR decomposition of  $\mathbf{L}$ , and the ALS method

guarantees that  $\mathcal{R}(\mathbf{L})$  is the left dominant subspace of  $\mathbf{A}_{(n)}$ . In many applications the orthogonal basis suffices, but in case the singular vectors are required, one can obtain them from the orthogonal basis by using, e.g., a low-overhead randomized approach [26]. Specifically, to calculate the singular vectors, line 5 of Algorithm 4 can be replaced with the following steps:

$$\begin{aligned} \mathbf{U}, \mathbf{\Sigma}, \mathbf{V}^T &\leftarrow \text{matrix SVD of } \hat{\mathbf{Q}}^T \mathbf{A}_{(n)}, \\ \mathbf{U}^{(n)} &\leftarrow \hat{\mathbf{Q}}\mathbf{U}. \end{aligned}$$

The ALS improved *st*-HOSVD algorithm, referred to as *st*-HOSVD-ALS, can be analogously derived, as presented in Algorithm 5.

---

**Algorithm 5** *st*-HOSVD-ALS

---

**Input:** Tensor  $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ , truncation  $(R_1, R_2, \dots, R_N)$

**Output:** Low multilinear rank approximation  $\hat{\mathcal{A}} \approx \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \dots \times_N \mathbf{U}^{(N)}$

- 1: Select an order of  $\{1, 2, \dots, N\}$ , i.e.,  $\{i_1, i_2, \dots, i_N\}$
  - 2:  $\mathcal{B} \leftarrow \mathcal{A}$
  - 3: **for all**  $n \in \{i_1, i_2, \dots, i_N\}$  **do**
  - 4:    $\mathbf{B}_{(n)} \leftarrow \mathcal{B}$  in matrix format
  - 5:    $\mathbf{L}, \mathbf{R} \leftarrow \text{ALS}(\mathbf{B}_{(n)}, R_n)$
  - 6:    $\hat{\mathbf{Q}}, \hat{\mathbf{R}} \leftarrow \text{reduced QR decomposition of } \mathbf{L}$
  - 7:    $\mathbf{U}^{(n)} \leftarrow \hat{\mathbf{Q}}$
  - 8:    $\mathbf{B}_{(n)} \leftarrow \hat{\mathbf{R}}\mathbf{R}^T$
  - 9:    $\mathcal{B} \leftarrow \mathcal{B}_{(n)}$  in tensor format
  - 10: **end for**
  - 11:  $\mathcal{G}_{(i_N)} \leftarrow \mathcal{B}_{(i_N)}$
  - 12:  $\mathcal{G} \leftarrow \mathcal{G}_{(i_N)}$  in tensor format
- 

The difference between Algorithm 4 and 5 is whether or not to store  $\mathbf{R}$  and  $\hat{\mathbf{R}}$ , the right factor matrices of the ALS method and reduced QR decomposition, respectively. Storing them will help reduce the overall computational cost when updating tensor  $\mathcal{B}$ , and core tensor  $\mathcal{G}$  can be calculated with the last factor matrix simultaneously in Algorithm 5. Apart from the computational cost of ALS in the *t*-HOSVD-ALS algorithm, calculating the core tensor  $\mathcal{G}$  is also critical, especially for higher-order tensors.

It is worth mentioning that the matricizations of  $\mathcal{A}$  and  $\mathcal{B}$  are not necessary if a row-wise update rule is used in *t*-HOSVD-ALS and *st*-HOSVD-ALS algorithms. Taking *t*-HOSVD-ALS as an example, in line 3 of Algorithm 4 we calculate the factor matrix  $\mathbf{U}^{(n)}$  by using an ALS method to obtain the rank- $R_n$  approximation of  $\mathbf{A}_{(n)}$ . The key computation is solving a multi-side least squares problem with the right-hand side  $\mathbf{A}_{(n)}$  or  $\mathbf{A}_{(n)}^T$ . And the multi-side least squares problem is equivalent to a series of independent least squares problems whose right-hand sides are the columns of  $\mathbf{A}_{(n)}$  or  $\mathbf{A}_{(n)}^T$ , i.e., the mode- $n$  fiber or slice of tensor  $\mathcal{A}$ . These independent least squares problems can be solved in parallel, and each least squares problem only requires a fiber or slice of tensor  $\mathcal{A}$ , therefore the explicit matricization in line 2 of Algorithm 4 can be naturally avoided.

Compared with *t*-HOSVD and *st*-HOSVD, the proposed algorithms exhibit several advantages. First, the redundant computations of the singular vectors are totally avoided, thus the overall cost of the algorithm can be substantially reduced. Second, the convergence of the ALS procedure is

controllable by adjusting the convergence tolerance. This is helpful considering the fact that  $t$ -HOSVD and  $st$ -HOSVD are quasi-optimal, and are often used as the initial guess for other iterative algorithms such as HOOI. Third, the algorithms are free of intermediate data explosion since the least square problems can be solved without any intermediate matrices.

An added benefit of the proposed  $t$ -HOSVD-ALS and  $st$ -HOSVD-ALS algorithms is that the solution of the multi-side least squares problems is intrinsically parallelizable. By using the ALS method, each row of the factor matrix  $\mathbf{L}$  or  $\mathbf{R}$  can be independently updated. Therefore, one can distribute the computation of the rows over multiple computing units. Since the workload for each row is almost identical, a simple static load distribution strategy suffices. All other operations in the algorithms, such as the matrix-matrix multiplication, the QR reduction and the small-scale matrix inversion, can also be easily parallelized by calling vendor-supplied highly optimized linear algebra libraries.

#### 4 Computational Cost and Error Analysis

In the proposed  $t$ -HOSVD-ALS and  $st$ -HOSVD-ALS algorithms, the performance of the ALS iteration depends on several factors, such as the initial guess and the convergence criterion. Based on the convergence property and the convergence condition of the ALS method, we suggest to set the initial guess  $\mathbf{L}_0$  as follows.

1. Generate a random matrix  $\mathbf{S}$ , whose entries are uniform distributions on interval  $[0, 1]$ .
2. Compute the reduced QR decomposition  $\mathbf{A}_{(n)}\mathbf{S} = \mathbf{Q}\mathbf{R}$ .
3. Let  $\mathbf{Q}$  be the initial guess, i.e.,  $\mathbf{L}_0 = \mathbf{Q}$ .

In this way, it is assured that  $\mathcal{R}(\mathbf{L}_0)$  is a subspace of  $\mathcal{R}(\mathbf{A}_{(n)})$ , which is closer to the left dominant subspace of  $\mathbf{A}_{(n)}$  than a random initial guess. Also, step 3 makes sure that the initial guess is properly normalized. We remark here that unlike techniques utilized in randomized algorithms that require  $\mathbf{S}$  to satisfy some special properties [26, 43, 44], the suggested initialization approach only requires  $\mathbf{S}$  to be dense and full rank.

The stopping condition of the ALS iteration can be set to

$$|\|\mathbf{A}_{(n)} - \mathbf{L}_k\mathbf{R}_k^T\|_F - \|\mathbf{A}_{(n)} - \mathbf{U}_1\mathbf{U}_1^T\mathbf{A}_{(n)}\|_F| \leq \eta\|\mathbf{A}\|_F, \quad (4.1)$$

where  $\mathcal{R}(\mathbf{U}_1)$  is the left dominant subspace of  $\mathbf{A}_{(n)}$ , and  $\eta$  is an accuracy tolerance parameter. In practice, however,  $\mathbf{U}_1$  is often not available. We therefore advise to replace (4.1) by

$$|\|\mathbf{A}_{(n)} - \mathbf{L}_k\mathbf{R}_k^T\|_F - \|\mathbf{A}_{(n)} - \mathbf{L}_{k+1}\mathbf{R}_{k+1}^T\|_F| \leq \eta\|\mathbf{A}\|_F \quad (4.2)$$

as the stop criterion, which means that the relative approximation error almost no longer decreases, implying that  $\{\mathbf{L}_k, \mathbf{R}_k\}$  has converged to the critical point of optimization problem  $\min_{\mathbf{L}, \mathbf{R}} \|\mathbf{A}_{(n)} - \mathbf{L}\mathbf{R}^T\|_F$ .

Next, we will discuss truncation  $R_n$  and how to select the tolerance parameter  $\eta$  by error analysis. To analyze the approximation error of ALS-based algorithms, we first recall a useful lemma.

**Lemma 2** [63] *Let  $\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times R_n}$ ,  $n \in \{1, 2, \dots, N\}$  be a sequence of column orthogonal matrices, calculated via the  $t$ -HOSVD or  $st$ -HOSVD algorithm, and suppose that  $\hat{\mathbf{A}} = \mathbf{A} \times_1 (\mathbf{U}^{(1)}\mathbf{U}^{(1)T}) \times_2 (\mathbf{U}^{(2)}\mathbf{U}^{(2)T}) \dots \times_N (\mathbf{U}^{(N)}\mathbf{U}^{(N)T})$  is an approximation of  $\mathbf{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ . Then*

$$\|\hat{\mathbf{A}} - \mathbf{A}\|_F^2 \leq \sum_{n=1}^N \gamma_n \leq N\|\mathbf{A} - \mathbf{A}_{opt}\|_F^2, \quad (4.3)$$

where  $\gamma_n = \sum_{r=R_n+1}^{I_n} (\sigma_r^{(n)})^2$ , and  $\mathcal{A}_{\text{opt}}$  is the optimal solution of problem (1.1).

It is worth noting that although estimation (4.3) ignores the computation error, it is still useful in practice. By Lemma 2, the error analysis of our algorithms is described in Theorem 3, with proof given in Appendix C.

**Theorem 3** *If the stop criterion of ALS is set to (4.1), then the approximation errors of Algorithm 4 and 5 are bounded by*

$$\frac{\|\hat{\mathcal{A}} - \mathcal{A}\|_F}{\|\mathcal{A}\|_F} \leq \sqrt{\sum_{n=1}^N (\eta_n^2 + \frac{\gamma_n}{\|\mathcal{A}\|_F^2})} \leq \sqrt{N}(\eta + \frac{\|\mathcal{A} - \mathcal{A}_{\text{opt}}\|_F}{\|\mathcal{A}\|_F}), \quad (4.4)$$

where  $\eta = \max_{n \in \{1, 2, \dots, N\}} \eta_n$ .

We remark that although in practice (4.1) is replaced by (4.2), numerical tests indicate that the main result (4.4) still holds. From (4.4), we advice to choose the tolerance parameter  $\eta_n$  such that the dominant term in the right hand side of (4.4) is  $\gamma_n/\|\mathcal{A}\|_F^2$  or  $\|\mathcal{A} - \mathcal{A}_{\text{opt}}\|_F/\|\mathcal{A}\|_F$ . Furthermore, if truncation  $R_n$  is selected appropriately, both  $\gamma_n$  and  $\eta_n$  will be small, and the ALS will converge very fast since  $\sigma_{R_n+1}/\sigma_{R_n} \ll 1$ . On the other hand, less suitable truncation  $R_n$  represents larger  $\gamma_n$  and therefore larger  $\eta_n$ , which in turn reduces the required number of ALS iterations.

Also of interest to us is the overall costs of the proposed algorithms. We analyze cases related to both general higher-order tensor  $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$  with truncation  $(R_1, R_2, \dots, R_N)$  and cubic tensor  $\mathcal{A} \in \mathbb{R}^{I \times I \times \dots \times I}$  with truncation  $(R, R, \dots, R)$ . The analysis results are shown in Table 1, where  $\text{iter}_n$  is the number of ALS iterations for mode  $n$ . The proposed ALS-based algorithms, the computational costs rely greatly on  $\text{iter}_n$ , which depends on the initial guess, the truncation  $R_n$  and the accuracy requirement. Our numerical result will reveal later that  $\text{iter}_n$  is usually far smaller than  $R_n$ , which is consistent with previous studies of the ALS method for matrix computation [60].

Table 1: Computational cost of different  $t$ - and  $st$ -HOSVD algorithms

Algorithms		$\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$	$\mathcal{A} \in \mathbb{R}^{I \times I \times \dots \times I}$
$t$ -HOSVD	ALS-based	$\mathcal{O}(\sum_{n=1}^N (R_n I_{1:N}) \text{iter}_n)$	$\mathcal{O}(\sum_{n=1}^N R I^N \text{iter}_n)$
	EIG-based	$\mathcal{O}(\sum_{n=1}^N (I_n I_{1:N} + R_n I_n^2))$	$\mathcal{O}(N I^{N+1} + N R I^2)$
	SVD-based	$\mathcal{O}(\sum_{n=1}^N R_n I_{1:N})$	$\mathcal{O}(N R I^N)$
$st$ -HOSVD	ALS-based	$\mathcal{O}(\sum_{n=1}^N (R_{1:n} I_{n:N}) \text{iter}_n)$	$\mathcal{O}(\sum_{n=1}^N (R^n I^{N-n+1}) \text{iter}_n)$
	EIG-based	$\mathcal{O}(\sum_{n=1}^N (R_{1:n-1} I_n I_{n:N} + R_n I_n^2))$	$\mathcal{O}(\sum_{n=1}^N R^{n-1} I^{N-n+2} + N R I^2)$
	SVD-based	$\mathcal{O}(\sum_{n=1}^N R_{1:n} I_{n:N})$	$\mathcal{O}(\sum_{n=1}^N R^n I^{N-n+1})$

For comparison purpose, we also list in the table the complexities of the original  $t$ - and  $st$ -HOSVD algorithms, including the ones based on matrix SVD and those based on the eigen-decomposition of the Gram matrix. From the table we can make the following observations.

- The EIG-based algorithms are usually most costly as compared to the SVD and ALS-based ones, especially when the truncation size is much smaller than the dimension length.
- It is hard to tell theoretically whether SVD-based algorithms are more or less costly than the proposed the ALS-based algorithms. We will examine their cost through numerical experiments in the next section.

## 5 Numerical Experiments

In this section, we will compare the proposed ALS-based algorithms with the original  $t$ -HOSVD and  $st$ -HOSVD algorithms by several numerical experiments related to both synthetic and real-world tensors. The implementation of the original  $t$ - and  $st$ -HOSVD algorithms includes `mlsvd` from Tensorlab [65] and `hosvd` from Tensor Toolbox [5]. In particular, `mlsvd` utilizes matrix SVD and `hosvd` employs eigen-decomposition of Gram matrix, for computing the factor matrices, therefore we denote them as  $t(st)$ -HOSVD-SVD and  $t(st)$ -HOSVD-EIG, respectively. To examine the numerical behaviors of these algorithms, we carry out most of the experiments in MATLAB R2019b on a computer equipped with an Intel Xeon Gold 6240 CPU of 2.60 GHz. And to study the parallel performance of the proposed algorithms, we implement the algorithms in C++ and run them on a workstation equipped with an Intel Xeon Gold 6154 CPU of 3.00 GHz. Unless mentioned otherwise, the tolerance parameter is set to  $\eta = 10^{-4}$ , and the maximum number of ALS iterations is limited to 50 in all tests.

### 5.1 Reconstruction of random low-rank tensors with noise

In the first set of experiments we examine the performance of the original  $t$ - and  $st$ -HOSVD algorithms, as well as the proposed ALS-based ones for the reconstruction of random low-rank tensors with noise. The tests are designed following the work of [72, 62]. Specifically, the input tensor is randomly generated as

$$\hat{\mathcal{A}} = \mathcal{A} + \delta \mathcal{E},$$

where the elements of  $\mathcal{E}$  follow the standard Gaussian distribution, and the noise level is controlled by  $\delta = 10^{-4}$ . The base tensor  $\mathcal{A}$  is constructed by two models, i.e., CP model and Tucker model. And to measure the reconstruction error we use  $\|\mathcal{B} - \mathcal{A}\|_F / \|\mathcal{A}\|_F$ , where  $\mathcal{B}$  is the low multilinear rank approximation of  $\hat{\mathcal{A}}$ .

#### 5.1.1 CP model

First, we use the CP model to construct the base tensor  $\mathcal{A} \in \mathbb{R}^{I \times J \times K}$  as follows:

$$\mathcal{A} = \lambda_1 \cdot \mathbf{a}_1 \circ \mathbf{b}_1 \circ \mathbf{c}_1 + \lambda_2 \cdot \mathbf{a}_2 \circ \mathbf{b}_2 \circ \mathbf{c}_2 + \cdots + \lambda_R \cdot \mathbf{a}_R \circ \mathbf{b}_R \circ \mathbf{c}_R,$$

where  $\mathbf{a}_r \in \mathbb{R}^I$ ,  $\mathbf{b}_r \in \mathbb{R}^J$ ,  $\mathbf{c}_r \in \mathbb{R}^K$  are randomly generated normalized vectors, and coefficients  $\lambda_r \in [5, 10]$  for all  $r \in \{1, 2, \dots, R\}$ .

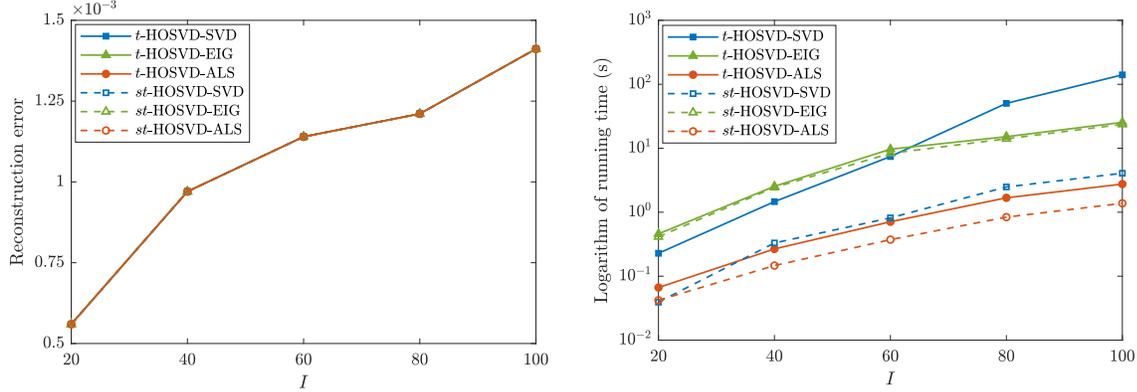


Fig. 1: Reconstruction errors and running time of various low multilinear rank approximation algorithms for reconstructing random noisy low CP-rank tensors with gradually increased size.

In the experiments, we set the tensor size to be  $J = I$  and  $K = 100I$  and gradually increase  $I$  from 20 to 100 with step 20. The truncation is set to  $(R, R, R)$ , where  $R = 0.2I$ . We carry out the tests for 20 times and draw the averaged reconstruction errors and running time in Fig. 1. From the figure, it is observed that there is almost no difference in reconstruction error among all tested algorithms, indicating that the proposed ALS-based methods can maintain the accuracy of the original ones. In terms of the running time,  $t$ -HOSVD-ALS is  $3.4\times \sim 50.9\times$  faster than  $t$ -HOSVD-SVD, and  $6.9\times \sim 13.6\times$  faster than  $t$ -HOSVD-EIG, respectively. For  $st$ -HOSVD, the speedup of  $st$ -HOSVD-ALS is  $1.0\times \sim 3.0\times$  and  $9.8\times \sim 22.3\times$ , as compared to  $st$ -HOSVD-SVD and  $st$ -HOSVD-EIG, respectively. We remark that the original algorithms behaves very differently in  $t$ -HOSVD and  $st$ -HOSVD. In particular, changing from  $t$ -HOSVD-EIG to  $st$ -HOSVD-EIG leads to little performance improvement. This is because the efficiency of eigen-decomposition of the Gram matrix is strongly dependent on the size of the third mode of the input tensor, and the sequentially updated algorithm is not able to help in this case.

### 5.1.2 Tucker model

Then we consider the base tensor  $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times I_3 \times I_4}$  constructed through the Tucker model in the following way:

$$\mathcal{A} = \mathcal{G} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \times_3 \mathbf{U}^{(3)} \times_4 \mathbf{U}^{(4)},$$

where  $\mathcal{G} \in \mathbb{R}^{R_1 \times R_2 \times R_3 \times R_4}$  is randomly generated core tensor whose elements follow the uniform distribution on interval  $[5, 10]$ , and  $\mathbf{U}^{(i)} \in \mathbb{R}^{I_i \times R_i}$ ,  $i = 1, 2, 3, 4$  are column orthogonal factor matrices.

In the experiments, we set  $I_2 = I_3 = I_4 = 100$  and gradually increase  $I_1 = I$  from 1,000 to 5,000 with step 1,000. The truncation is set to be  $(R, 10, 10, 10)$ , where  $R = 0.01I$ . We again carry out the tests for 20 times, and draw the averaged reconstruction errors and running time in Fig. 2. From the figure, we observe that the reconstruction errors of all tested algorithms are nearly identical. For  $t$ -HOSVD, the performance of  $t$ -HOSVD-ALS is close to that of  $t$ -HOSVD-EIG, and they are both  $10.3\times$  faster than  $t$ -HOSVD-SVD. For  $st$ -HOSVD, the speedup of  $st$ -HOSVD-ALS is  $26.5\times$  and  $3.5\times$ , as compared to  $st$ -HOSVD-SVD and  $st$ -HOSVD-EIG, respectively.

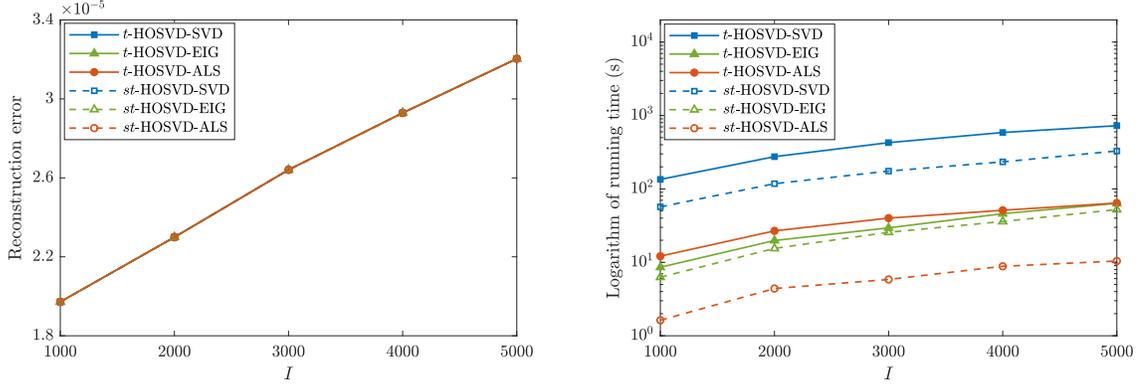


Fig. 2: Reconstruction errors and running time of various low multilinear rank approximation algorithms for reconstructing random noisy low multilinear-rank tensors with gradually increased size.

## 5.2 Classification of handwritten digits

The second set of experiments is designed for testing the capability of the original  $t$ - and  $st$ -HOSVD algorithms and the proposed ALS-based ones on handwritten digits classification. It was studied that low multilinear rank approximation can be applied to compress the training data of images so that the core tensor can be utilized for image classification on the test data [10]. In the tests, we use the MNIST database [39, 40] of handwritten digits. We transfer the training dataset into a fourth-order tensor  $\mathcal{A} \in \mathbb{R}^{28 \times 28 \times 5000 \times 10}$ , where the first- and second-mode are the texel modes, the third-mode corresponds to training images, and the fourth-mode represents image categories. In the tests, the truncation is fixed to be (8, 8, 142, 10), which is close to the setting of the reference work [10]. We measure the classification accuracy of a classification algorithm as the percentage of the test images that is correctly classified.

Table 2: Approximation error, classification accuracy and training time of various low multilinear rank approximation algorithms for handwritten digits classification of the MNIST database.

Algorithms	$t$ -HOSVD			$st$ -HOSVD		
	SVD	EIG	ALS	SVD	EIG	ALS
Approximation error	0.4330	0.4330	0.4335	0.4288	0.4288	0.4291
Classification accuracy (%)	95.45	95.45	95.45	95.36	95.36	95.35
Training time (s)	18.97	5.38	3.30	4.07	5.12	1.53

We run the test for 20 times and record the averaged approximation error, classification accuracy and running time of the tested algorithms in Table 2. From the table it can be seen that the approximation errors and classification accuracies of all tested algorithms are almost indistin-

guishable, which again validates the accuracy of the proposed methods. Table 2 also shows that the ALS-based methods are fastest in terms of training time. Specifically, the ALS-based approaches are on average  $1.6\times \sim 5.7\times$  and  $2.7\times \sim 3.3\times$  faster than the original  $t$ -HOSVD and  $st$ -HOSVD algorithms, respectively.

### 5.3 Compression of tensors arising from fluid dynamics simulations

The purpose of this set of experiments is to examine the performance of different low multilinear rank approximation algorithms for compressing tensors generated from the simulation results of a lid-driven cavity flow, which is a standard benchmark for incompressible fluid dynamics [11]. The simulation is done in a square domain of length 1 m with the speed of the top plate setting to 1 m/s and all other boundaries no slip. The kinematic viscosity is  $\nu = 1.0 \times 10^{-4} \text{ m}^2/\text{s}$ , and the fluid properties is assumed to be laminar. We use the OpenFOAM software package [2] to conduct the simulation on a uniform grid with 100 grid cells in each direction. The simulation is run with time step  $\Delta t = 1.0 \times 10^{-4} \text{ s}$  and terminated at  $t = 1.0 \text{ s}$ . We record the magnitude of velocity at every time step. The simulation results of the lid-driven cavity flow are stored in a third-order tensor of size  $100 \times 100 \times 10000$ . To test the tensor approximation algorithms, we fix the truncation to  $(20, 20, 20)$ , corresponding to a compression ratio of 12,500 : 1.

First, we examine whether  $(R_1, R_2, R_3) = (20, 20, 20)$  is a suitable truncation for the input tensor. Since the dimensions of the first two modes are both 100, we consider  $R_1 = R_2 = 20$  is a relatively proper choice. And  $R_3 = 20$  is remained to test for the third dimension of length 10,000. Therefore, we perform a test to see how the approximation error varies as  $R_3$  changes gradually from 10 to 100. The test results are shown in Fig. 3, from which it can be observed that  $R_3 = 20$  is also proper for the third dimension.

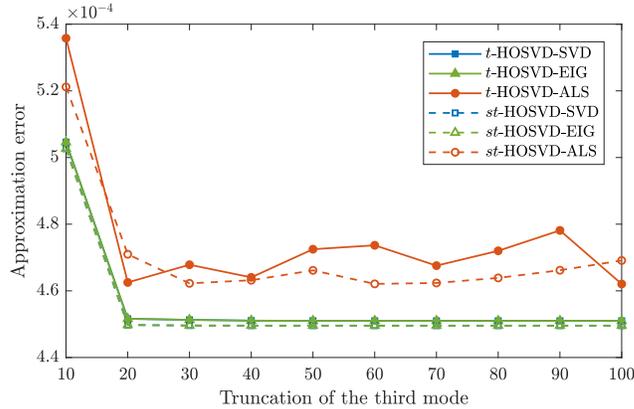


Fig. 3: Approximation errors of  $t$ -HOSVD and  $st$ -HOSVD with gradually increased  $R_3$ .

Next, we study the efficiency of tested algorithms under different accuracy requirements. We run the test for 20 times with tolerance parameter  $\eta$  adjusted to different values and record the averaged relative residual and running time for each value of  $\eta$ ; the test results are listed in Table 3.

Table 3: Relative residuals and running time of various low multilinear rank approximation algorithms for compressing tensors arising from fluid dynamics simulations with different tolerance parameters.

$t$ -HOSVD	SVD	EIG	ALS		
			$\eta = 10^{-2}$	$\eta = 10^{-4}$	$\eta = 10^{-6}$
Relative residual ( $\times 10^{-4}$ )	4.5161	4.5161	15.7422	4.6240	4.5225
Running time (s)	118.35	23.61	1.41 + 0.96	2.21 + 0.96	3.964 + 0.96
$st$ -HOSVD	SVD	EIG	ALS		
			$\eta = 10^{-2}$	$\eta = 10^{-4}$	$\eta = 10^{-6}$
Relative residual ( $\times 10^{-4}$ )	4.4976	4.4976	14.2730	4.7139	4.5044
Running time (s)	3.44	21.70	0.73	1.16	1.82

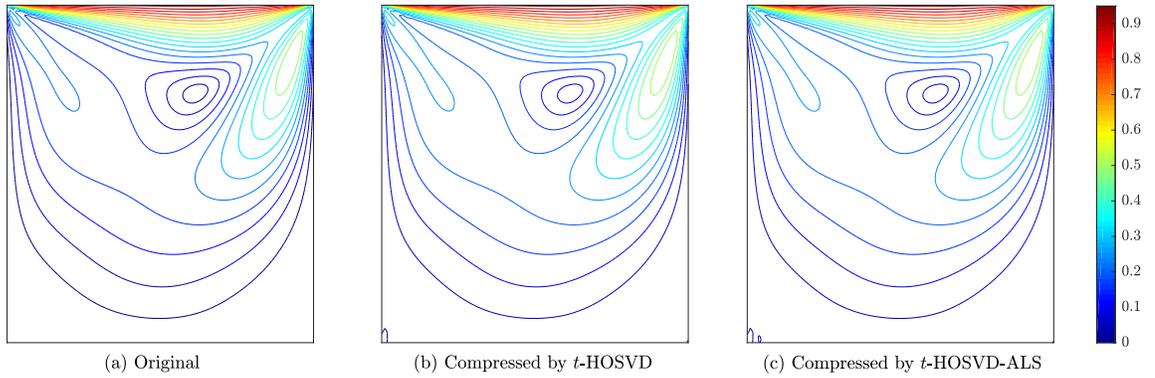


Fig. 4: The original and compressed data at  $t = 0.5$  s for compressing tensors arising from fluid dynamics simulations with tolerance parameter  $\eta = 10^{-4}$ .

Also provided in the table is the extra cost of computing the singular vectors for  $t$ -HOSVD-ALS, if requested. From the table we have the following observations.

- The relative residuals and running time of the original  $t$ - and  $st$ -HOSVD algorithms are independent of the change of the tolerance parameter  $\eta$ . This is due to the usage of Krylov subspace method for computing matrix truncated SVD or eigen-decomposition.
- For the ALS-based methods, the relative residuals and the running time both depend on  $\eta$ . With  $\eta$  decreased, the relative residuals are reduced to a similar level that the original algorithms can attain but more running time is required.
- The ALS-based algorithms are the fastest in all tests. It can achieve  $6.1 \times \sim 45.6 \times$  speedup for  $t$ -HOSVD and  $2.3 \times \sim 18.4 \times$  speedup for  $st$ -HOSVD, respectively.
- The overhead of computing the singular vectors for  $t$ -HOSVD-ALS is independent of the ALS tolerance and is relatively low.

It seems from the tests that despite the excellent performance of the proposed ALS-based methods, the strong dependency between the sustained performance and the tolerance parameter  $\eta$  could

eventually lead to poor performance if  $\eta$  is very small. In practice the  $t$ - or  $st$ -HOSVD algorithm is often used as the initial guess of a supposedly more accurate iterative method such as HOOI. In this case it is not necessary to use a very tight tolerance parameter. To examine whether  $\eta = 10^{-4}$  is a suitable choice for the ALS-based algorithms in the same tests, we draw in Fig. 4 the contours of the original and compressed velocity data at  $t = 0.5$  s. It clearly shows that when  $\eta = 10^{-4}$ , the compressed results are consistent with each other with very little discrepancy. In fact, the measured maximum differences between the original data and the compressed data obtained by  $t$ -HOSVD and that by  $t$ -HOSVD-ALS are around  $1.58 \times 10^{-3}$  and  $1.33 \times 10^{-3}$ , respectively, which are very small considering that the compression ratio is over five orders of magnitude.

Table 4: A comparison of HOOI results for compressing tensors arising from fluid dynamics simulations with initial solutions provided by various low multilinear rank approximation algorithms.

Algorithm	Relative residual		Number of HOOI iterations
	Initial	Final	
$t$ -HOSVD-SVD	$4.5161 \times 10^{-4}$	$4.4807 \times 10^{-4}$	9.0
$t$ -HOSVD-EIG	$4.5161 \times 10^{-4}$	$4.4807 \times 10^{-4}$	7.0
$t$ -HOSVD-ALS	$4.6260 \times 10^{-4}$	$4.4807 \times 10^{-4}$	6.0
$st$ -HOSVD-SVD	$4.4976 \times 10^{-4}$	$4.4807 \times 10^{-4}$	5.0
$st$ -HOSVD-EIG	$4.4976 \times 10^{-4}$	$4.4807 \times 10^{-4}$	6.0
$st$ -HOSVD-ALS	$4.5044 \times 10^{-4}$	$4.4807 \times 10^{-4}$	7.0

To further investigate the applicability of the compressed results, we use the computed low multilinear rank approximation with tolerance parameter  $\eta = 10^{-4}$  as the initial guess of the HOOI method with stopping criterion  $10^{-12}$ . The HOOI method is obtained from the `tucker_als` function of the Tensor Toolbox v3.1 [5]. The test results are presented in Table 4, in which we list the relative residuals with the  $t$ - and  $st$ -HOSVD provided initial guesses, the final relative residuals of HOOI, and the number of HOOI iterations, all averaged on 20 independent runs. From the table we can see that although the initial residual provided by the ALS-based algorithms are slightly larger than those provided by the original  $t$ - and  $st$ -HOSVD methods, same final residuals can be achieved after HOOI iterations nevertheless. And more importantly, the required numbers of HOOI iterations are insensitive to which specific  $t$ - and  $st$ -HOSVD algorithms, original or not, are used as shown in the tests. In other words, the proposed ALS-based methods are able to deliver similar results as the original ones when applying in HOOI, even when the tolerance parameter is relatively loose.

#### 5.4 Parallel performance

An advantage of the proposed ALS-based methods is that they are easy to parallelize. To examine the parallel performance, in this experiment, we implement the  $t$ -HOSVD-ALS and  $st$ -HOSVD-ALS algorithms in C++ with OpenMP multi-threading parallelization [3]. The involved linear algebra operations are available with parallelization from the Intel MKL [1, 67] and the open-source ARMADILLO [51, 52] libraries. In the test, the input tensor is generated by `ttensor` in

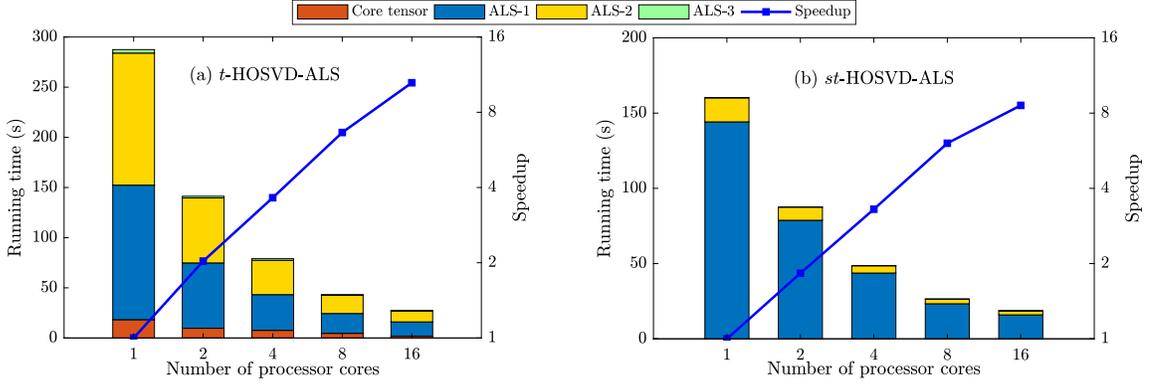


Fig. 5: The running time and speedup of the ALS-based truncated HOSVD algorithms for compressing low multilinear rank tensor on a parallel computer.

the Tensor Toolbox, whose size is  $5000 \times 5000 \times 50$  and the multilinear rank is  $(500, 500, 5)$ . We break down the running time into different portions, including the ALS iterations along the three dimensions (ALS- $i$ ,  $i \in \{1, 2, 3\}$ ) and the calculation of the core tensor. The test results are drawn in Fig. 5. Also shown in the figures is the parallel scalability of the proposed algorithms.

From the figure, it can be seen that the proposed  $t$ -HOSVD-ALS and  $st$ -HOSVD-ALS can both scale well. In particular,  $t$ -HOSVD-ALS and  $st$ -HOSVD-ALS can achieve speedups of  $10.50\times$  and  $8.59\times$  as the number of processor cores is increased from 1 to 16, respectively. Moreover, in both algorithms the ALS iterations, as the major costs, are accelerated efficiently with the increased number of processor cores. Another observation is that there is a slight drop of parallel efficiency when using 16 processor cores, which is caused by the fact that the factor matrix  $U^{(n)}$  is usually tall and skinny due to the low multilinear rank structure of tensor. Overall, the test results demonstrate that the proposed ALS-based algorithms are parallelization friendly and have good potential to scale further on larger high-performance computers.

## 6 Conclusions

In this paper, we proposed a class of ALS-based algorithms for efficiently calculating the low multilinear rank approximation of tensors. Compared with the original  $t$ -HOSVD and  $st$ -HOSVD algorithms, the proposed algorithms are superior in several ways. First, by eliminating the redundant computations of the singular vectors, the overall costs of the algorithms are substantially reduced. Second, the proposed algorithms are more flexible with adjustable convergence tolerance, which is especially useful when the algorithms are used to generate initial solutions for iterative methods such as HOOI. Third, the proposed algorithms are free of the notorious data explosion issue due to the fact that the ALS procedure does not explicitly require the intermediate matrices. And fourth, the ALS-based approaches are parallelization friendly on high-performance computers. Theoretical analysis shows that the ALS iteration in the proposed algorithms is  $q$ -linear convergent with a relatively wide convergence region. Numerical experiments with both synthetic and real-world tensor data demonstrate that proposed ALS-based algorithms can substantially reduce the total cost of low multilinear rank approximation and are highly parallelizable.

Possible future works could include applying of the proposed ALS-based algorithms to more applications, among which we are especially interested in large-scale scientific computing. It would also be of interest to study randomization techniques to further improve the performance of the proposed algorithms, considering the fact that solving multiple least squares problems with different right-hand sides is the major cost. Some of the ideas presented in this work, such as the utilization of ALS for solving the intermediate low rank approximation problem, might be extended to other tensor decomposition models such as tensor-train (TT) and hierarchical Tucker (HT) decompositions.

## Acknowledgments

The authors would like to thank the anonymous reviewers for their valuable comments and suggestions that have greatly improved the quality of this paper.

## Declarations

**Funding** This study was funded in part by Guangdong Key R&D Project (#2019B121204008), Beijing Natural Science Foundation (#JQ18001) and Beijing Academy of Artificial Intelligence.

**Conflicts of Interest** The authors declare that they have no conflict of interest.

**Availability of Data and Material** The datasets generated and analyzed during the current study are available from the corresponding author on reasonable request.

**Code Availability** The code used in the current study is available from the corresponding author on reasonable request.

## References

1. Intel Math Kernel Library Reference Manual [EQ/OL]. URL <http://developer.intel.com>
2. OpenFoam 2018, Version 6. URL <https://openfoam.org>
3. OpenMP application program interface, Version 4.0, OpenMP Architecture Review Board (July, 2013). URL <http://www.openmp.org>
4. Austin, W., Ballard, G., Kolda, T.G.: Parallel tensor compression for large-scale scientific data. In: IEEE International Parallel and Distributed Processing Symposium, pp. 912–922. IEEE (2016)
5. Bader, B.W., Kolda, T.G., et.al.: MATLAB Tensor Toolbox Version 3.1. Available online (2019). URL <https://www.tensortoolbox.org>
6. Baglama, J., Reichel, L.: Augmented implicitly restarted Lanczos bidiagonalization methods. SIAM J. Sci. Comput. **27**(1), 19–42 (2005)
7. Beckmann, C., Smith, S.: Tensorial extensions of independent component analysis for multisubject fMRI analysis. Neuroimage **25**(1), 294–311 (2005)
8. Beylkin, G., Mohlenkamp, M.J.: Numerical operator calculus in higher dimensions. Proc. Natl. Acad. Sci. **99**, 10246–10251 (2002)
9. Bro, R.: Review on multiway analysis in chemistry—2000–2005. Crit. Rev. Anal. Chem. **36**, 279–293 (2006)
10. B.Savas, Eldén, L.: Handwritten digit classification using higher order singular value decomposition. Pattern Recognition **40**, 993–1003 (2007)
11. Burggraf, R.: Analytical and numerical studies of the structure of steady separated flows. J. Fluid Mechanics **24**(1), 113–151 (1966)
12. Carroll, J.D., Chang, J.J.: Analysis of individual differences in multidimensional scaling via an  $N$ -way generalization of “Eckart-Young” decomposition. Psychometrika **35**, 283–319 (1970)

13. Charalampous, K., Gasteratos, A.: A tensor-based deep learning framework. *Image and Vision Computing* **32**(11), 916–929 (2014)
14. Cullum, J., Willoughby, R., Lake, M.: A Lanczos algorithm for computing singular values and vectors of large matrices. *SIAM J. Sci. Stat. Comput.* **4**(2), 197–215 (1983)
15. De Lathauwer, L., De Moor, B., Vandewalle, J.: A multilinear singular value decomposition. *SIAM J. Matrix Anal. Appl.* **21**(4), 1253–1278 (2000)
16. De Lathauwer, L., De Moor, B., Vandewalle, J.: On the best rank-1 and rank- $(r_1, r_2, \dots, r_N)$  approximation of higher-order tensors. *SIAM J. Matrix Anal. Appl.* **21**, 1324–1342 (2000)
17. De Lathauwer, L., Vandewalle, J.: Dimensionality reduction in higher-order signal processing and rank- $(r_1, r_2, \dots, r_N)$  reduction in multilinear algebra. *Linear Algebra Appl.* **391**, 31–55 (2004)
18. De Leeuw, J., Young, F., Takane, Y.: Additive structure in qualitative data: An alternating least squares method with optimal scaling features. *Psychometrika* **41**, 471–503 (1976)
19. 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)
20. Ding, W., Wei, Y.: Solving multi-linear systems with  $\mathcal{M}$ -tensors. *J. Sci. Comput.* **68**, 689–715 (2016)
21. Elden, L., Savas, B.: A Newton-Grassmann method for computing the best multilinear rank- $(r_1, r_2, r_3)$  approximation of a tensor. *SIAM J. Matrix Anal. Appl.* **31**(2), 248–271 (2009)
22. Golub, G.H., Van Loan, C.F.: *Matrix Computations* Ed. 4th. Johns Hopkins University Press, Baltimore (2008)
23. Grasedyck, L.: Hierarchical singular value decomposition of tensors. *SIAM J. Matrix Anal. Appl.* **31**(4), 2029–2054 (2010)
24. Hackbusch, W.: Numerical tensor calculus. *Acta Numerica* **23**, 651–742 (2014)
25. Hackbusch, W., Kühn, S.: A new scheme for the tensor representation. *J. Fourier Anal. Appl.* **15**, 706–722 (2009)
26. Halko, N., Martinsson, P.G., Tropp, J.A.: Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Rev.* **53**(2), 217–288 (2011)
27. Henrion, R.: Body diagonalization of core matrices in three-way principal components analysis: Theoretical bounds and simulation. *J. Chemometrics* **7**, 477–494 (1993)
28. Hitchcock, F.L.: Multiple invariants and generalized rank of a  $p$ -way matrix or tensor. *J. Math. Phys.* **7**(1-4), 39–79 (1928)
29. Holtz, S., Rohwedder, T., Schneider, R.: The alternating linear scheme for tensor optimization in the tensor train format. *SIAM J. Sci. Comput.* **34**(2), A683–A713 (2012)
30. Ishteva, M., Absil, P.A., Van Huffel, S., De Lathauwer, L.: Best low multilinear rank approximation of higher-order tensors, based on the Riemannian trust-region scheme. *SIAM J. Matrix Anal. Appl.* **31**(1), 115–135 (2011)
31. Ishteva, M., De Lathauwer, L., Absil, P.A., Huffel, S.: Dimensionality reduction for higher-order tensors: Algorithms and applications. *Int. J. Pure Appl. Math.* **42**(3), 337–343 (2012)
32. Ishteva, M., De Lathauwer, L., Absil, P.A., Huffel, S.V.: Differential-geometric Newton method for the best rank- $(R_1, R_2, R_3)$  approximation of tensors. *Numer. Algor.* **51**, 179–194 (2009)
33. Jiang, J., Wu, H., Li, Y., Yu, R.: Three-way data resolution by alternating slice-wise diagonalization (ASD) method. *J. Chemometrics* **14**, 15–36 (2000)
34. Khoromskij, B.N.: Tensors-structured numerical methods in scientific computing: Survey on recent advances. *Chemometrics and Intelligent Laboratory Systems* **110**(1), 1–19 (2012)
35. Kiers, H.A.L.: Towards a standardized notation and terminology in multiway analysis. *J. Chemometrics* **14**, 105–122 (2000)
36. Kolda, T.G., Bader, B.W.: Tensor decompositions and applications. *SIAM Rev.* **51**(3), 455–500 (2009)
37. Kroonenberg, P.M.: *Applied Multiway Data Analysis*. John Wiley & Sons, Inc., New Jersey (2008)
38. Kroonenberg, P.M., De Leeuw, J.: Principal component analysis of three-mode data by means of alternating least squares algorithms. *Psychometrika* **45**, 69–97 (1980)
39. LeCun, Y., C.Cortes, C.Burges: THE MNIST DATABASE of handwritten digits. URL <http://yann.lecun.com/exdb/mnist/>
40. LeCun, Y., L.Bottou, Y.Bengio, P.Haffner: Gradient-based learning applied to document recognition. *Proceedings of the IEEE* **84**(11), 2278–2324 (1998)
41. Levin, J.: Three-mode factor analysis. Ph.D. thesis, University of Illinois, Urbana-Champaign (1963)
42. Li, J., Battaglini, C., Perros, I., Sun, J., Vuduc, R.: An input-adaptive and in-place approach to dense tensor-times-matrix multiply. In: *Proceedings of the International Conference for High Performance Computing*, 76, pp. 1–12. ACM (2015)
43. Mahoney, M.W.: Randomized algorithms for matrices and data. *Foundations and Trends in Machine Learning* **3**(2), 123–224 (2011)
44. Mahoney, M.W., Drineas, P.: RandNLA: Randomized numerical linear algebra. *Communications of the ACM* **59**(6), 80–90 (2016)

45. Novikov, A., Podoprikin, D., Osokin, A., Vetrov, D.: Tensorizing neural networks. In: Annual Conference on Neural Information Processing Systems, vol. 1, pp. 442–450. ACM (2015)
46. Oh, J., Shin, K., Papalexakis, E., Faloutsos, C., Yu, H.: S-HOT: Scalable high-Order Tucker decomposition. In: Proceedings of the Tenth ACM International Conference on Web Search and Data Mining, pp. 761–770. ACM (2017)
47. Oh, S., Park, N., Jang, J., Sael, L., Kang, U.: High-performance tucker factorization on heterogeneous platforms. *IEEE Transactions on Parallel and Distributed Systems* **30**(10), 2237–2248 (2019)
48. Oseledets, I.V., Tyrtshnikov, E.E.: Breaking the curse of dimensionality, or how to use SVD in many dimensions. *SIAM J. Sci. Comput.* **31**(5), 3744–3759 (2009)
49. Oseledets, I.V.: Tensor-train decomposition. *SIAM J. Sci. Comput.* **33**(5), 2295–2317 (2011)
50. Rajbhandari, S., Nikam, A., Lai, P.W., Stock, K., Krishnamoorthy, S., Sadayappan, P.: A communication-optimal framework for contracting distributed tensors. In: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, pp. 375–386. IEEE (2014)
51. Sanderson, C., Curtin, R.: Armadillo: A template-based C++ library for linear algebra. *J. Open Source Software* **1**(2), 26 (2016)
52. Sanderson, C., Curtin, R.: A user-friendly hybrid sparse matrix class in C++. In: International Conference on Mathematical Software, Lecture Notes in Computer Science (LNCS), vol. 10931, pp. 422–430. Springer (2018)
53. Savas, B., Lim, L.H.: Quasi-Newton methods on Grassmannians and multilinear approximations of tensors. *SIAM J. Sci. Comput.* **32**(6), 3352–3393 (2009)
54. Schatz, M.: Distributed tensor computations: Formalizing distributions, redistributions, and algorithm derivations. Ph.D. thesis, University of Texas, Austin (2015)
55. Shashua, A., Levin, A.: Linear image coding for regression and classification using the tensor-rank principle. In: IEEE Conference on Computer Vision and Pattern Recognition, pp. 42–49. IEEE (2001)
56. Sidiropoulos, N.D., De Lathauwer, L., Xiao, F., Huang, K.J., Papalexakis, E.E., Faloutsos, C.: Tensor decomposition for signal processing and machine learning. *IEEE Trans. on Sig. Proc.* **65**(13), 3551–3582 (2017)
57. Smith, S., Ravindran, N., Sidiropoulos, N., Karpypis, G.: Efficient and parallel sparse tensor matrix multiplication. In: IEEE International Parallel and Distributed Processing Symposium, pp. 61–70. IEEE (2015)
58. Sorensen, D.: Implicit application of polynomial filters in a  $k$ -step Arnoldi method. *SIAM J. Matrix Anal. Appl.* **13**(1), 357–385 (1992)
59. Stewart, G.: A Krylov-Schur algorithm for large eigenproblems. *SIAM J. Matrix Anal. Appl.* **23**(3), 601–614 (2001)
60. Szlam, A., Tulloch, A., Tygert, M.: Accurate low-rank approximations via a few iterations of alternating least squares. *SIAM J. Matrix Anal. Appl.* **38**(2), 425–433 (2017)
61. Tucker, L.R.: Some mathematical notes on three-mode factor analysis. *Psychometrika* **31**, 279–311 (1966)
62. Vannieuwenhoven, N., Vandebril, R., Meerbergen, K.: On the truncated multilinear singular value decomposition. Technical Report TW589, Department of Computer Science, Katholieke Universiteit Leuven, Leuven, Belgium (2011)
63. Vannieuwenhoven, N., Vandebril, R., Meerbergen, K.: A new truncation strategy for the higher-order singular value decomposition. *SIAM J. Sci. Comput.* **34**(2), A1027–A1052 (2012)
64. Vasilescu, M.A.O., Terzopoulos, D.: Multilinear image analysis for facial recognition. In: Object recognition supported by user interaction for service robots, vol. 2, pp. 511–514. IEEE (2002)
65. Vervliet, N., Debals, O., Sorber, L., M. Van Barel, L. De Lathauwer: Tensorlab 3.0. Available online (2016). URL <https://www.tensorlab.net>
66. Vlasi, D., Brand, M., Pfister, H., Popovic, J.: Face transfer with multilinear models. In: ACM SIGGRAPH 2005 Papers, vol. 24, pp. 426–433. ACM (2005)
67. Wang, E., Zhang, Q., Shen, B., Zhang, G., Lu, X., Wu, Q., Wang, Y.: “Intel Math Kernel Library,” in High-Performance Computing on the *Intel<sup>®</sup> Xeon Phi*. Springer pp. 167–188 (2014)
68. Wang, S.G., Wu, M.X., Jia, Z.Z.: Matrix Inequality Ed. 2th. Science Press, Beijing (2006)
69. Watkins, D.: The QR algorithm revisited. *SIAM Rev.* **50**(1), 133–145 (2008)
70. Xu, Y.: On the convergence of higher-order orthogonal iteration. *Linear and Multilinear Algebra* **66**(11), 2247–2265 (2018)
71. Young, F.W., Takane, Y., De Leeuw, J.: The principal components of mixed measurement level multivariate data: An alternating least squares method with optimal scaling features. *Psychometrika* **43**, 279–281 (1978)
72. Zhang, T., Golub, G.H.: Rank-one approximation to high order tensors. *SIAM J. Matrix Anal. Appl.* **23**(2), 534–550 (2001)

## A Proof of Theorem 1

Based on the assumption of Theorem 1,  $\mathbf{L}_k$  is nonsingular and  $\mathbf{L}_k^T \mathbf{L}_k$  is positive definite. Thus, the iterative form of Algorithm 3 is

$$\begin{aligned}\mathbf{R}_k &= \mathbf{A}^T \mathbf{L}_k (\mathbf{L}_k^T \mathbf{L}_k)^{-1}, \\ \mathbf{L}_{k+1} &= \mathbf{A} \mathbf{R}_k (\mathbf{R}_k^T \mathbf{R}_k)^{-1},\end{aligned}\tag{A.1}$$

i.e.,

$$\mathbf{L}_{k+1} = \mathbf{A} \mathbf{A}^T \mathbf{L}_k (\mathbf{L}_k^T \mathbf{A} \mathbf{A}^T \mathbf{L}_k)^{-1} (\mathbf{L}_k^T \mathbf{L}_k).\tag{A.2}$$

Suppose that the full SVD of  $\mathbf{A}$  is  $\mathbf{A} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T$ , where

$$\mathbf{U} = [\mathbf{U}_1, \mathbf{U}_2], \quad \mathbf{V} = [\mathbf{V}_1, \mathbf{V}_2], \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_1 & 0 \\ 0 & \boldsymbol{\Sigma}_2 \end{pmatrix}.$$

Then from (A.2), we have

$$\mathbf{U}^T \mathbf{L}_{k+1} = \mathbf{U}^T \mathbf{A} \mathbf{V} \mathbf{V}^T \mathbf{A}^T \mathbf{U} \mathbf{U}^T \mathbf{L}_k (\mathbf{L}_k^T \mathbf{U} \mathbf{U}^T \mathbf{A} \mathbf{V} \mathbf{V}^T \mathbf{A}^T \mathbf{U} \mathbf{U}^T \mathbf{L}_k)^{-1} (\mathbf{L}_k^T \mathbf{U} \mathbf{U}^T \mathbf{L}_k),\tag{A.3}$$

which can be rewritten into block form

$$\begin{pmatrix} \mathbf{L}_{k+1}^{(1)} \\ \mathbf{L}_{k+1}^{(2)} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\Sigma}_1^2 & 0 \\ 0 & \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T \end{pmatrix} \begin{pmatrix} \mathbf{L}_k^{(1)} \\ \mathbf{L}_k^{(2)} \end{pmatrix} (\mathbf{L}_k^{(1)T} \boldsymbol{\Sigma}_1^2 \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T \mathbf{L}_k^{(2)})^{-1} (\mathbf{L}_k^{(1)T} \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \mathbf{L}_k^{(2)}).$$

It then follows that

$$\begin{aligned}\mathbf{L}_{k+1}^{(1)} &= \boldsymbol{\Sigma}_1^2 \mathbf{L}_k^{(1)} (\mathbf{L}_k^{(1)T} \boldsymbol{\Sigma}_1^2 \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T \mathbf{L}_k^{(2)})^{-1} (\mathbf{L}_k^{(1)T} \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \mathbf{L}_k^{(2)}), \\ \mathbf{L}_{k+1}^{(2)} &= \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T \mathbf{L}_k^{(2)} (\mathbf{L}_k^{(1)T} \boldsymbol{\Sigma}_1^2 \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T \mathbf{L}_k^{(2)})^{-1} (\mathbf{L}_k^{(1)T} \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \mathbf{L}_k^{(2)}).\end{aligned}\tag{A.4}$$

Furthermore,

$$\begin{aligned}& (\mathbf{L}_k^{(1)T} \boldsymbol{\Sigma}_1^2 \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T \mathbf{L}_k^{(2)})^{-1} \\ &= (\mathbf{I} + (\mathbf{L}_k^{(1)T} \boldsymbol{\Sigma}_1^2 \mathbf{L}_k^{(1)})^{-1} (\mathbf{L}_k^{(2)T} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T \mathbf{L}_k^{(2)}))^{-1} (\mathbf{L}_k^{(1)T} \boldsymbol{\Sigma}_1^2 \mathbf{L}_k^{(1)})^{-1} \\ &= (\mathbf{I} + \sum_{n=1}^{\infty} (-1)^n ((\mathbf{L}_k^{(1)T} \boldsymbol{\Sigma}_1^2 \mathbf{L}_k^{(1)})^{-1} (\mathbf{L}_k^{(2)T} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T \mathbf{L}_k^{(2)}))^n) (\mathbf{L}_k^{(1)T} \boldsymbol{\Sigma}_1^2 \mathbf{L}_k^{(1)})^{-1}.\end{aligned}\tag{A.5}$$

Here we suppose that the distance between  $\mathcal{R}(\mathbf{L}_k)$  and  $\mathcal{R}(\mathbf{U}_2)$  is small enough, therefore can be denoted as  $\delta_k$ , which only depends on  $\|\mathbf{L}_k^{(2)}\|_2$  (i.e., there exist two constants  $\alpha, \beta > 0$  such that  $\alpha \delta_k \leq \|\mathbf{L}_k^{(2)}\|_2 \leq \beta \delta_k$ ). We can then obtain the lower bound of the distance between  $\mathcal{R}(\mathbf{L}_k)$  and  $\mathcal{R}(\mathbf{U}_1)$ , which is  $\sqrt{1 - \delta_k^2}$ , and

$$\|\mathbf{L}_k^{(1)}\|_2 \leq C_1, \quad \|(\mathbf{L}_k^{(1)})^{-1}\|_2 \leq \frac{C_2}{\sqrt{1 - \delta_k^2}},\tag{A.6}$$

where  $C_1, C_2$  are constants independent on  $k$  and  $\delta_k$ . From (A.5) and (A.6), there exists a constant  $C$  that is only dependent on  $C_1, C_2$  so that the following inequality holds.

$$(\mathbf{L}_k^{(1)T} \boldsymbol{\Sigma}_1^2 \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T \mathbf{L}_k^{(2)})^{-1} \leq (\mathbf{L}_k^{(1)T} \boldsymbol{\Sigma}_1^2 \mathbf{L}_k^{(1)})^{-1} + \frac{C \delta_k^2}{(1 - 2\delta_k^2)^2},\tag{A.7}$$

Further, from (A.7) and (A.4), assume that  $\sigma_r > \sigma_{r+1}$ , we have

$$\mathbf{L}_{k+1}^{(2)} \leq \frac{\boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T}{\sigma_r^2} \mathbf{L}_k^{(2)} (\mathbf{L}_k^{(1)T} \frac{\boldsymbol{\Sigma}_1^2}{\sigma_r^2} \mathbf{L}_k^{(1)})^{-1} (\mathbf{L}_k^{(1)T} \mathbf{L}_k^{(1)}) + \hat{C} \left( \frac{\delta_k^2}{(1 - 2\delta_k^2)^2} + \frac{\delta_k^2}{1 - \delta_k^2} + \frac{\delta_k^4}{(1 - 2\delta_k^2)^2} \right),$$

where  $\tilde{C}$  is a constant. Clearly,

$$0 < \mathbf{L}_k^{(1)T} \mathbf{L}_k^{(1)} \leq \mathbf{L}_k^{(1)T} \frac{\boldsymbol{\Sigma}_1^2}{\sigma_r^2} \mathbf{L}_k^{(1)},$$

and by Lemma 1, we have

$$\|(\mathbf{L}_k^{(1)T} \frac{\boldsymbol{\Sigma}_1^2}{\sigma_r^2} \mathbf{L}_k^{(1)})^{-1} (\mathbf{L}_k^{(1)T} \mathbf{L}_k^{(1)})\|_2 \leq 1.$$

Since  $\delta_k$  is small enough, we obtain

$$\|\mathbf{L}_{k+1}^{(2)}\|_2 \leq \frac{\sigma_{r+1}^2}{\sigma_r^2} \|\mathbf{L}_k^{(2)}\|_2 + \tilde{C} \delta_k^2 \leq \frac{\sigma_{r+1}^2}{\sigma_r^2} \|\mathbf{L}_k^{(2)}\|_2 + \frac{\tilde{C}}{\alpha^2} \|\mathbf{L}_k^{(2)}\|_2^2, \quad (\text{A.8})$$

where  $\alpha, \tilde{C}$  do not depend on  $k$  and  $\|\mathbf{L}_k^{(2)}\|_2$ .

Denote

$$q = \frac{\sigma_{r+1}^2}{\sigma_r^2} + \frac{\tilde{C}}{\alpha^2} \|\mathbf{L}_0^{(2)}\|_2.$$

Since we assume that  $\mathcal{R}(\mathbf{L}_0)$  is close to  $\mathcal{R}(\mathbf{U}_1)$  enough,  $\|\mathbf{U}_2^T \mathbf{L}_0\|_2$  is sufficiently small, i.e.,  $\|\mathbf{L}_0^{(2)}\|_2 = o(1)$ . In other words, we assume that  $q < 1$ . From (A.8), we have

$$\|\mathbf{L}_{k+1}^{(2)}\|_2 \leq q \|\mathbf{L}_k^{(2)}\|_2$$

for all  $k$ , which leads to

$$\lim_{k \rightarrow +\infty} \|\mathbf{L}_k^{(2)}\|_2 \rightarrow 0.$$

Combining with the assumption of  $\mathbf{L}_k$ , it is verified that  $\mathcal{R}(\mathbf{L}_k)$  is orthogonal to  $\mathcal{R}(\mathbf{U}_2)$  with  $k \rightarrow +\infty$ . Since the orthogonal complement space of  $\mathcal{R}(\mathbf{U}_2)$  is unique, we have

$$\mathcal{R}(\mathbf{L}_k) = \mathcal{R}(\mathbf{U}_1), \quad k \rightarrow +\infty,$$

where  $\mathcal{R}(\mathbf{U}_1)$  is the dominant subspace of  $\mathbf{A}$ . In other words, we have

$$\lim_{k \rightarrow +\infty} \|\mathbf{L}_k \mathbf{L}_k^\dagger - \mathbf{U}_1 \mathbf{U}_1^T\|_2 = 0,$$

where  $\mathbf{L}_k^\dagger$  is the pseudo-inverse of  $\mathbf{L}_k$ . Further, from the iterative form of the ALS method, we have

$$\mathbf{R}_k = \mathbf{A}^T (\mathbf{L}_k^\dagger)^T,$$

thus

$$\mathbf{L}_k \mathbf{R}_k^T = \mathbf{L}_k \mathbf{L}_k^\dagger \mathbf{A} \rightarrow \mathbf{U}_1 \mathbf{U}_1^T \mathbf{A}, \quad k \rightarrow +\infty.$$

And since the Frobenius norm  $\|\cdot\|_F$  is continuous,

$$\lim_{k \rightarrow +\infty} \|\mathbf{A} - \mathbf{L}_k \mathbf{R}_k^T\|_F = \|\mathbf{A} - \mathbf{U}_1 \mathbf{U}_1^T \mathbf{A}\|_F.$$

Since  $\mathbf{U}_1 \mathbf{U}_1^T \mathbf{A}$  is the exact solution of low rank approximation of  $\mathbf{A}$ , the convergence of the ALS method is proved.

From (A.8), we further confirm the  $q$ -linear convergence of the ALS method, with approximate convergence ratio  $\sigma_{r+1}^2/\sigma_r^2$ .  $\square$

## B Proof of Theorem 2

The assumption of Theorem 1 implies that  $\mathbf{L}_k$  is nonsingular at every iteration  $k$ . We assume that

$$\mathbf{L}_k^{(1)T} \boldsymbol{\Sigma}_1^2 \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T \mathbf{L}_k^{(2)}$$

is positive definite. Let  $\varepsilon$  be a positive number such that  $\sigma_r > \sigma_{r+1} - \varepsilon$ . By (A.4), we know

$$\begin{aligned}\mathbf{L}_{k+1}^{(1)} &= \boldsymbol{\Sigma}_1^2 \mathbf{L}_k^{(1)} (\mathbf{L}_k^{(1)T} \boldsymbol{\Sigma}_1^2 \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T \mathbf{L}_k^{(2)})^{-1} (\mathbf{L}_k^{(1)T} \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \mathbf{L}_k^{(2)}), \\ \mathbf{L}_{k+1}^{(2)} &= \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T \mathbf{L}_k^{(2)} (\mathbf{L}_k^{(1)T} \boldsymbol{\Sigma}_1^2 \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T \mathbf{L}_k^{(2)})^{-1} (\mathbf{L}_k^{(1)T} \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \mathbf{L}_k^{(2)}),\end{aligned}\tag{B.1}$$

which means

$$\begin{aligned}\mathbf{L}_{k+1}^{(1)} &= \frac{\boldsymbol{\Sigma}_1^2}{(\sigma_r - \varepsilon)^2} \mathbf{L}_k^{(1)} (\mathbf{L}_k^{(1)T} \frac{\boldsymbol{\Sigma}_1^2}{(\sigma_r - \varepsilon)^2} \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \frac{\boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T}{(\sigma_r - \varepsilon)^2} \mathbf{L}_k^{(2)})^{-1} (\mathbf{L}_k^{(1)T} \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \mathbf{L}_k^{(2)}), \\ \mathbf{L}_{k+1}^{(2)} &= \frac{\boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T}{(\sigma_r - \varepsilon)^2} \mathbf{L}_k^{(2)} (\mathbf{L}_k^{(1)T} \frac{\boldsymbol{\Sigma}_1^2}{(\sigma_r - \varepsilon)^2} \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \frac{\boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T}{(\sigma_r - \varepsilon)^2} \mathbf{L}_k^{(2)})^{-1} (\mathbf{L}_k^{(1)T} \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \mathbf{L}_k^{(2)}).\end{aligned}$$

Clearly it holds that

$$\|\mathbf{L}_{k+1}^{(2)}\|_2 \leq \frac{\sigma_{r+1}^2}{(\sigma_r - \varepsilon)^2} \|\mathbf{L}_k^{(2)}\|_2\tag{B.2}$$

under the condition that

$$\mathbf{L}_k^{(1)T} \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \mathbf{L}_k^{(2)} \leq \mathbf{L}_k^{(1)T} \frac{\boldsymbol{\Sigma}_1^2}{(\sigma_r - \varepsilon)^2} \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \frac{\boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T}{(\sigma_r - \varepsilon)^2} \mathbf{L}_k^{(2)}.\tag{B.3}$$

If  $\mathbf{L}_k^{(1)}$  is nonsingular, then (B.3) implies

$$(\mathbf{L}_k^{(2)} (\mathbf{L}_k^{(1)})^{-1})^T (\mathbf{I} - \frac{\boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T}{(\sigma_r - \varepsilon)^2}) (\mathbf{L}_k^{(2)} (\mathbf{L}_k^{(1)})^{-1}) \leq \frac{\boldsymbol{\Sigma}_1^2}{(\sigma_r - \varepsilon)^2} - \mathbf{I}.\tag{B.4}$$

It follows to see that

$$\|\mathbf{L}_k^{(2)} (\mathbf{L}_k^{(1)})^{-1}\|_2 \leq \sqrt{\frac{\sigma_r^2 - (\sigma_r - \varepsilon)^2}{(\sigma_r - \varepsilon)^2 - \sigma_{min}^2}}\tag{B.5}$$

is a sufficient condition of (B.4).

Next we will prove that if the initial guess  $\mathbf{L}_0$  satisfies condition (B.5), then  $\mathbf{L}_k^{(1)}$  is nonsingular and (B.2) is satisfied at every iteration  $k$ .

Provided that  $\mathbf{L}_0$  satisfies (B.5), we obtain

$$\|\mathbf{L}_1^{(2)}\|_2 \leq \frac{\sigma_{r+1}^2}{(\sigma_r - \varepsilon)^2} \|\mathbf{L}_0^{(2)}\|_2.\tag{B.6}$$

And according to the proof of Theorem 1, we know  $\mathbf{L}_1^{(1)}$  is also nonsingular, which implies

$$\mathbf{L}_1^{(1)T} \boldsymbol{\Sigma}_1^2 \mathbf{L}_1^{(1)} + \mathbf{L}_1^{(2)T} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T \mathbf{L}_1^{(2)}$$

is positive definite. Then by (B.1), we have

$$\|\mathbf{L}_1^{(2)} (\mathbf{L}_1^{(1)})^{-1}\|_2 \leq \|\boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T \mathbf{L}_0^{(2)} (\mathbf{L}_0^{(1)})^{-1} \boldsymbol{\Sigma}_1^{-2}\|_2 \leq \frac{\sigma_{r+1}^2}{\sigma_r^2} \|\mathbf{L}_0^{(2)} (\mathbf{L}_0^{(1)})^{-1}\|_2 \leq \sqrt{\frac{\sigma_r^2 - (\sigma_r - \varepsilon)^2}{(\sigma_r - \varepsilon)^2 - \sigma_{min}^2}}.\tag{B.7}$$

Analogously, we can prove that for every iteration  $k$ ,  $\mathbf{L}_k^{(1)}$  is nonsingular, i.e.,  $\mathbf{L}_k^{(1)T} \boldsymbol{\Sigma}_1^2 \mathbf{L}_k^{(1)} + \mathbf{L}_k^{(2)T} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_2^T \mathbf{L}_k^{(2)}$  is positive definite, and (B.5) is satisfied. Since (B.5) is a sufficient condition of (B.4) and (B.3), inequality (B.2) is true at every iteration  $k$ , which implies that

$$\lim_{k \rightarrow 0} \|\mathbf{L}_k \mathbf{L}_k^\dagger - \mathbf{U}_1 \mathbf{U}_1^T\|_2 = 0.$$

The rest part of the proof is analogous to the proof of Theorem 1, which is omitted for brevity.  $\square$

### C Proof of Theorem 3

In Algorithm 4,  $\mathbf{U}^{(n)}$  is obtained from the rank- $R_n$  approximation of  $\mathbf{A}_{(n)}$ , which is done in an iterative manner and allows a tolerance parameter  $\eta_n$ . Therefore, we have

$$\|\mathbf{A}_{(n)} - \mathbf{L}^* \mathbf{R}^{*T}\|_F^2 \leq \eta_n^2 \|\mathbf{A}\|_F^2 + \gamma_n, \quad (\text{C.1})$$

where  $\mathbf{L}^*$  and  $\mathbf{R}^*$  are the same as in Algorithm 4. Note that  $\mathbf{R}$  is updated by solving a multi-side least squares problem

$$\min_{\mathbf{R}} \|\mathbf{A}_{(n)} - \mathbf{L} \mathbf{R}^T\|_F,$$

whose exact solution is  $\mathbf{R} = \mathbf{A}_{(n)}^T (\mathbf{L}^\dagger)^T$ . Thus

$$\|\mathbf{A}_{(n)} - \mathbf{L}^* \mathbf{R}^{*T}\|_F = \|\mathbf{A}_{(n)} - \mathbf{L}^* \mathbf{L}^{*\dagger} \mathbf{A}_{(n)}\|_F, \quad (\text{C.2})$$

where  $\mathbf{L}^* \mathbf{L}^{*\dagger}$  represents an orthogonal projection on subspace  $\mathcal{R}(\mathbf{L}^*)$ . Consequently, by (C.1), (C.2) and (4.3), we have

$$\|\hat{\mathbf{A}} - \mathbf{A}\|_F^2 \leq \sum_{n=1}^N (\eta_n^2 \|\mathbf{A}_{(n)}\|_F^2 + \gamma_n),$$

which means

$$\frac{\|\hat{\mathbf{A}} - \mathbf{A}\|_F^2}{\|\mathbf{A}\|_F^2} \leq \sum_{n=1}^N \left( \eta_n^2 + \frac{\gamma_n}{\|\mathbf{A}\|_F^2} \right).$$

Combining with (4.3), we obtain (4.4).

The error analysis of Algorithm 5 can be analogously done. □