# Orthogonal Nonnegative Tucker Decomposition

Junjun Pan, Michael K. Ng, Ye Liu, Xiongjun Zhang, Hong Yan*†‡§

### Abstract

In this paper, we study the nonnegative tensor data and propose an orthogonal nonnegative Tucker decomposition (ONTD). We discuss some properties of ONTD and develop a convex relaxation algorithm of the augmented Lagrangian function to solve the optimization problem. The convergence of the algorithm is given. We employ ONTD on the image data sets from the real world applications including face recognition, image representation, hyperspectral unmixing. Numerical results are shown to illustrate the effectiveness of the proposed algorithm.

**Keywords.** nonnegative tensor, Tucker decomposition, image processing

## 1   Introduction

Given a nonnegative matrix $\mathbf{A} \in \mathbb{R}_+^{m \times n}$ and integer $r$, nonnegative matrix factorization (NMF) is the problem of searching for basic matrix $\mathbf{U} \in \mathbb{R}_+^{m \times r}$ and coefficient matrix $\mathbf{V} \in \mathbb{R}_+^{r \times n}$ such that $\mathbf{A} \approx \mathbf{UV}$. In many data analysis problems, the columns of $\mathbf{A}$ are corresponding to data points, for instance, images of pixel intensities. NMF has been successfully applied into many fields including image processing, text data mining and so on. It has been demonstrated that NMF is a powerful technique for dimension reduction. Compared to other well-known method, like singular value decomposition or principal component analysis, NMF is able to give more interpretable results due to its combinations of nonnegative basic vectors.

In general, NMF is NP-hard and the solution is not unique. It is necessary to impose additional constraints on the factor matrix like orthogonality constraints. Precisely, given $\mathbf{A} \in \mathbb{R}_+^{m \times n}$, solve

$$\min_{\mathbf{U} \in \mathbb{R}_+^{m \times r}, \mathbf{V} \in \mathbb{R}_+^{r \times n}} \|\mathbf{A} - \mathbf{UV}\|_F^2, \quad \text{s.t.} \quad \mathbf{U}^T \mathbf{U} = \mathbf{I}.$$

We call the above problem orthogonal nonnegative matrix factorization(ONMF). The orthogonal constraints guarantee the uniqueness of the solution. There are many methods [1–3] and most are the multiplicative update algorithms derived from NMF. Until recently, Pan and Ng [4] investigated the properties of ONMF and present a new method called SN-ONMF for finding the factorization. They used the sparsity and nuclear norm optimization tucker to solve ONMF problem.

The orthogonality constraints make sense in many practical applications. In [3], the equivalence of ONMF problem and K-means clustering has been well discussed. In document classification, each entry $A(i,j)$ indicates the importance of word $i$ in document $j$. Each row of data matrix

---

*J. Pan is with Department of Mathematics and Operational Research Faculté polytechnique, Université de Mons, 7000, Belgium. (e-mail: Junjun.PAN@umons.ac.be).

†M. K. Ng, and Y. Liu are with the Department of Mathematics, The University of Hong Kong, Pokfulam, Hong Kong. (e-mail: mng@math.hku.hk; 16482549@life.hkbu.edu.hk).

‡X. Zhang is with the School of Mathematics and Statistics and Hubei Key Laboratory of Mathematical Sciences, Central China Normal University, Wuhan 430079, China. (e-mail: xjzhang@mail.ccnu.edu.cn).

§H. Yan is with the Department of Electronic Engineering, City University of Hong Kong, Kowloon, Hong Kong (e-mail:h.yan@cityu.edu.hk).

stands for a document, each column stands for a word. $\mathbf{U}$ is the document cluster indicator matrix in ONMF model, in other words, ONMF aims to find a document clustering space $\mathbf{U}$, and the coefficient matrix $\mathbf{V}$ can be obtained by projecting the data onto $\mathbf{U}$. In addition to understanding the cluster of documents, one may also require the cluster of words. Considering this problem, Ding et al proposed the following nonnegative tri-factor decomposition in [3],

$$\min_{\mathbf{U}\in\mathbb{R}_+^{m\times r_1},\mathbf{S}\in\mathbb{R}^{r_1\times r_2},\mathbf{V}\in\mathbb{R}_+^{r_2\times n}} \|\mathbf{A}-\mathbf{USV}\|_F^2, \quad \text{s.t.} \quad \mathbf{U}^T\mathbf{U}=\mathbf{I}, \quad \mathbf{V}^T\mathbf{V}=\mathbf{I}. \tag{1}$$

$\mathbf{U}$ provides row clusters and $\mathbf{V}$ provides column clusters. For coefficient matrix $\mathbf{S}$, each entry $s_{i,j}$ can be regarded as the connection weight between column cluster $i$ and row cluster $j$.

Nowadays, data that comes from many fields are more naturally represented as multidimensional data which refers to tensor, for example, video data, hyperspectral data, fMRI data and so on. In this paper, we generalize model (1) to tensor data, i.e., given a nonnegative tensor $\mathcal{A}\in\mathbb{R}_+^{I_1\times I_2\times\cdots\times I_d}$ and the integer rank $(J_1, J_2, \ldots, J_d)$, solve

$$\min \|\mathcal{A}-\mathcal{S}\times_1\mathbf{U}^{(1)}\times_2\mathbf{U}^{(2)}\cdots\times_d\mathbf{U}^{(d)}\|_F^2$$
$$\text{s.t. } \mathcal{S}\in\mathbb{R}_+^{J_1\times J_2\times\cdots\times J_d}, \ \mathbf{U}^{(n)}\in\mathbb{R}_+^{I_n\times J_n}, \ \mathbf{U}^{(n)T}\mathbf{U}^{(n)}=\mathbf{I}, \ n=1,2,\ldots,d. \tag{2}$$

where $\times_n$ denotes the mode-$n$ matrix product of a tensor defined by

$$(\mathcal{S}\times_n\mathbf{U}^{(n)})_{j_1\cdots j_{n-1}i_n j_{n+1}\cdots j_d} = \sum_{j_n=1}^{J_n} s_{j_1\cdots j_{n-1}j_n j_{n+1}\cdots j_d}u_{i_n,j_n}^{(n)}.$$

For simplicity, we call the above model orthogonal nonnegative Tucker decomposition (ONTD) model. $\mathbf{U}^{(n)}$ gives the clusters of the $n$-th dimension. Each entry $s_{j_1,j_2,\cdots,j_d}$ represents the joint connection weight of the corresponding cluster along dimensions from 1 to $d$. If some factor matrices are equal to identity matrix $\mathbf{I}$, we call the model partial ONTD model.

The ONTD model makes sense in several applications. In image classification, for some image sequences containing different illuminations, motions and subjects, ONTD model gives illumination clusters, motion clusters and subject clusters. We can also know the connection weight between them. For video data sets which contains different types of human actions, scenarios and different subjects, ONTD model helps us to know the clusters of actions, scenarios and subjects. If one only consider the clusters of actions and subjects, he can use partial ONTD model, i.e., set the scenarios factor matrix $\mathbf{U}$ to be identity matrix.

ONTD model not only helps to keep the inherent tensor structure but also well performs in data compression. As we all know that tensor data need huge storages due to high dimensionality. While in ONTD, only $(J_1\cdots J_d + \sum_{n=1}^d I_n J_n)$ memories are required for the tensor of $(I_1,\cdots,I_d)$. It would save a lot of memory compared with the original storage $(I_1 I_2\cdots I_d)$, especially when $(J_1,\cdots,J_d)$ is small.

ONTD model is related to multilinear singular value decomposition (HOSVD) which enforces the factor matrices in Tucker decomposition into orthogonal matrix and is well discussed in [5,6]. Higher-order orthogonal iteration (HOOI) in [5] is proposed to find the best rank $(J_1, J_2,\cdots,J_d)$ approximation of tensor $\mathcal{A}$. The difference between HOSVD and our model is the nonnegativity of factor matrices are imposed in ONTD model, which lead to easily interpret. In [7], Kim and Choi developed nonnegative Tucker decomposition (NTD) while they did not consider the orthogonality on the factor matrices. ONTD model takes advantages from both orthogonality and nonnegativity constraints. We can get the clustering information from the factor matrices $\{U^{(i)}\}_{i=1}^d$ and their joint connection weight from the core tensor $\mathcal{S}$ at the same time.

To solve the model (2), we first discuss the properties of nonnegative orthogonal Tucker format tensor, i.e., $\mathcal{A} = \mathcal{S} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \cdots \times_d \mathbf{U}^{(d)}$ where $\{\mathbf{U}^{(n)}\}_{n=1}^d$ are orthogonal nonnegative matrices. We utilize the properties of nonnegative orthonormal factor matrix and present a structured convex optimization algorithm. The convergence of the algorithm will be discussed and shown. The proposed ONTD method is then applied to face recognition, image representation and hyperspectral unmixing problems. Numerical results demonstrate a good performance of our model. We summarize the main contributions of this paper as follows.

1) We propose an orthogonal nonnegative Tucker decomposition model (ONTD) which projects the tensor objects into the tensor with small size for dimension reduction and preserves the structure information as well.

2) We present some properties of ONTD model, develop a structured convex optimization algorithm and analysis the convergence.

3) Numerical examples from the real world applications have been conducted to demonstrate the effectiveness of the proposed method. The results show that the ONTD outperforms existing methods such as PCA, NMF and NTD.

The rest of the paper is organized as follows. In Section 2, we present several properties of orthogonal nonnegative tensors. In Section 3 we propose an optimization model and present an algorithm. Meanwhile, the convergence is discussed. In Section 4, we apply the algorithm on real data sets from face recognition, image representation and hyperspectral unmixing problem. Their numerical results are shown. The concluding remarks are given in Section 5.

## 2  Properties of Nonnegative Orthogonal Tucker Tensor

Given a nonnegative orthogonal Tucker tensor $\mathcal{A} \in \mathbb{R}^{I_1 \times \cdots \times I_d}$,

$$\mathcal{A} = \mathcal{S} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \cdots \times_d \mathbf{U}^{(d)} \tag{3}$$

$\{\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times J_n}\}_{n=1}^d$ are orthogonal nonnegative matrices, $(J_1, J_2, \cdots, J_n)$ refers to as multilinear rank. We first introduce the following two properties of orthogonal nonnegative matrix that from [4].

**Lemma 1.** *For $n = 1, 2, \ldots, d$, each row of $\mathbf{U}^{(n)}$ has at most one nonzero element.*

**Lemma 2.** *For $n = 1, 2, \ldots, d$, $\mathbf{K}^{(n)} = \mathbf{U}^{(n)}\mathbf{U}^{(n)T}$ satisfies*

$$\mathbf{K}^{(n)T} = \mathbf{K}^{(n)} \text{ and } (\mathbf{K}^{(n)})^2 = \mathbf{K}^{(n)}, \text{ with } 0 \le k_{i,j}^{(n)} \le 1.$$

*In addition, the trace of $\mathbf{K}^{(n)}$ is equal to $J_n$, i.e., $tr(\mathbf{K}^{(n)}) = J_n$ and $\|\mathbf{K}^{(n)}\|_1 \le I_n$, where $I_n$ and $J_n$ are the dimensions of $\mathbf{U}^{(n)}$. Moreover, 1 is the $J_n$ repeated eigenvalues of $\mathbf{K}^{(n)}$ and the columns of $\mathbf{U}^{(n)}$ are the corresponding eigenvectors.*

From Lemma 1, we know $\mathbf{U}^{(n)}$ has the following structure:

$$\mathbf{U}^{(n)} = \Pi_r \begin{pmatrix} \mathbf{u}_1^{(n)} & 0 & \cdots & 0 \\ 0 & \mathbf{u}_2^{(n)} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{u}_{J_n}^{(n)} \end{pmatrix} \Pi_c, \quad \mathbf{u}_j^{(n)} = \begin{pmatrix} u_{1,j}^{(n)} \\ u_{2,j}^{(n)} \\ \vdots \\ u_{l_j,j}^{(n)} \end{pmatrix}, \tag{4}$$

where $j = 1, \cdots, J_n$, $n = 1, 2, \cdots, d$, and $\sum_j l_j = I_n$. $\Pi_c$ and $\Pi_r$ are permutation matrix. $\mathbf{U}^{(n)}$ can be regarded as class indicator matrix along $n$-th direction, $\mathbf{u}_j^{(n)}$ represents $j$-th class containing $l_j$ rows.

3

For tensor $\mathcal{A}$, its $n$-th unfolding $\mathbf{A}_{(n)} \in \mathbb{R}^{I_n \times (I_{n+1} \cdots I_N I_1 \cdots I_{n-1})}$, follow (3), we can say that the $I_n$ rows of $\mathbf{A}_{(n)}$ can be classified into $J_n$ classes. $\mathbf{U}^{(n)}$ is its corresponding class indicator matrix. For simplicity, denote $\mathbf{A}_{(n)}(t,:)$ as the $t$-th row of $\mathbf{A}_{(n)}$, and $\mathbf{A}_{(n)}(T_j,:)$ be the row set belongs to $j$-th class, with the row index set $T_j$ of cardinality $l_j$, $j \in \{1, \cdots J_n\}$.

Authors in [4] proved that the rows and columns of the same groups are proportional if matrix $\mathbf{A}$ is orthogonal decomposable, i.e., $\mathbf{A} = \mathbf{B}\mathbf{C}$ where $\mathbf{B}$ is orthogonal nonnegative matrix. Utilize this result, we give the similar property of $\mathcal{A}$ as follows.

**Property 1.** *Given a nonnegative orthogonal Tucker tensor $\mathcal{A}$, for $n$-th unfolding $\mathbf{A}_{(n)}$, $n = 1, 2, \ldots, d$, the rows belong to the same class are proportional.*

*Proof.* From Lemma 1, we know that $\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times J_n}$ has at most one nonzero element in each row, i.e., in the $i$-th row, $u_{i,j}^{(n)} = \begin{cases} u_{i,j^*}^{(n)}, & \text{if } j = j^*, \\ 0, & \text{otherwise}. \end{cases}$ From Tucker format, we get that

$$a_{i_1,i_2,\cdots,i_d} = \sum_{j_1,j_2,\cdots,j_d} s_{j_1,j_2,\cdots,j_d} u_{i_1,j_1}^{(1)} \cdots, u_{i_d,j_d}^{(d)} = s_{j_1^*,\cdots,j_d^*} u_{i_1,j_1^*}^{(1)} \cdots, u_{i_d,j_d^*}^{(d)} \tag{5}$$

Fix all the index but $i_n$,

$$\frac{a_{i_1,\cdots,i_{n-1},i_n^1,i_{n+1}\cdots i_d}}{a_{i_1,\cdots,i_{n-1},i_n^2,i_{n+1}\cdots i_d}} = \frac{s_{j_1^*,\cdots,j_n^{*1}\cdots j_d^*} u_{i_n^1 j_n^{*1}}^{(n)}}{s_{j_1^*,\cdots,j_n^{*2}\cdots j_d^*} u_{i_n^2 j_n^{*2}}^{(n)}},$$

if $j_n^{*1} = j_n^{*2}$, then,

$$\frac{a_{i_1,\cdots,i_{n-1},i_n^1,i_{n+1}\cdots i_d}}{a_{i_1,\cdots,i_{n-1},i_n^2,\cdots i_{n+1}i_d}} = \frac{u_{i_n^1 j_n^{*1}}^{(n)}}{u_{i_n^2 j_n^{*2}}^{(n)}}, \tag{6}$$

(6) implies that $i_n^1$-th and $i_n^2$-th row of $\mathbf{A}_{(n)}$ are in the same class, and proportional. The result follows. $\qquad\square$

**Property 2.** *If $\mathcal{A}$ is a nonnegative orthogonal Tucker tensor, then its factor matrices $\{\mathbf{U}^{(n)}\}_{n=1}^d$ can be given explicitly by $\{\mathbf{A}_{(n)}\}_{n=1}^d$.*

*Proof.* For $U^{(n)}$ defined in (4), the $j$-th column $\mathbf{u}_j^{(n)} = (\ u_{1,j}^{(n)} \ \ u_{2,j}^{(n)} \ \ \cdots \ \ u_{l_j,j}^{(n)}\ )^T$ can be constructed in the following way. Without loss of generality, let $u_{1,j}^{(n)} \neq 0$, from (6),

$$\frac{u_{t,j}^{(n)}}{u_{1,j}^{(n)}} = \frac{a_{i_1,\cdots,i_{n-1},t,i_{n+1},\cdots,i_d}}{a_{i_1,\cdots,i_{n-1},1,i_{n+1},\cdots,i_d}} \doteq \alpha_{t,j},$$

where $t \in \{1, 2, \cdots, l_j\}$, $j \in \{1, 2, \cdots, J_n\}$. Because of $\mathbf{U}^{(n)T}\mathbf{U}^{(n)} = I$, $\mathbf{u}_j^{(n)}$ can be constructed by letting

$$u_{t,j}^{(n)} = \frac{\alpha_{t,j}}{\alpha_j}, \qquad \alpha_j = \sqrt{\sum_{t=1}^{l_j} \alpha_{t,j}^2}.$$

$\qquad\square$

**Property 3.** *If $\mathcal{A}$ is a nonnegative orthogonal Tucker tensor, then for core tensor $\mathcal{S}$, the Frobenius norm of $j$-th row of $\mathbf{S}_{(n)}$ equals to that of $\mathbf{A}_{(n)}(T_j,:)$. More over, the Frobenius norm of $\mathcal{S}$ equals to that of $\mathcal{A}$.*

4

*Proof.* For $\mathbf{A}_{(n)}$ whose rows are classified into $J_n$ classes, let $\{a_{i_1,\cdots,i_{n-1},t,i_{n+1},\cdots,i_d}\}_{t=1}^{l_j}$ belong to the $j$-th class, from (5),

$$
\begin{aligned}
\|\mathbf{A}_{(n)}(T_j,:)\|_F^2 &= \sum_{i_1,\cdots,i_{n-1},t,i_{n+1},\cdots,i_d} a_{i_1,\cdots,i_{n-1},t,i_{n+1},\cdots,i_d}^2 \\
&= \sum_{i_1,\cdots,i_{n-1},t,i_{n+1},\cdots,i_d} (s_{j_1^*,\cdots,j_d^*} u_{i_1,j_1^*} \cdots u_{i_d,j_d^*})^2 \\
&= \sum_{i_1,\cdots,i_{n-1},i_{n+1},\cdots,i_d} \sum_t (s_{j_1^*,\cdots,j_d^*}^2 u_{t,j}^2) u_{i_1,j_1^*}^2 \cdots u_{i_{n-1},j_{n-1}^*}^2 u_{i_{n+1},j_{n+1}^*}^2 \cdots u_{i_d,j_d^*}^2 \\
&= \sum_{i_1,\cdots,i_{n-1},i_{n+1},\cdots,i_d} s_{j_1^*,\cdots,j_{n-1}^*,j,j_{n+1}^*,\cdots,j_d^*}^2 u_{i_1,j_1^*}^2 \cdots u_{i_{n-1},j_{n-1}^*}^2 u_{i_{n+1},j_{n+1}^*}^2 \cdots u_{i_d,j_d^*}^2 \quad (7) \\
&= \sum_{i_2,\cdots,i_{n-1},i_{n+1}\cdots i_d} \left(\sum_{i_1} s_{j_1^*,\cdots,j_{n-1}^*,j,j_{n+1}^*,\cdots,j_d^*}^2 u_{i_1,j_1^*}^2\right) \cdots u_{i_{n-1},j_{n-1}^*}^2 u_{i_{n+1},j_{n+1}^*}^2 \cdots u_{i_d,j_d^*}^2 \\
&= \sum_{i_2,\cdots,i_{n-1},i_{n+1}\cdots i_d} \left(\sum_{j_1^*} s_{j_1^*,\cdots,j_{n-1}^*,j,j_{n+1}^*,\cdots,j_d^*}^2\right) \cdots u_{i_{n-1},j_{n-1}^*}^2 u_{i_{n+1},j_{n+1}^*}^2 \cdots u_{i_d,j_d^*}^2 \quad (8) \\
&= \cdots = \sum_{j_1^*,\cdots,j_{n-1}^* j_{n+1}^* \cdots j_d^*} s_{j_1^*,\cdots,j_{n-1}^*,j,j_{n+1}^*,\cdots,j_d^*}^2 = \|\mathbf{S}_{(n)}(j,:)\|_F^2.
\end{aligned}
$$

Because of $\mathbf{U}^{(n)T}\mathbf{U}^{(n)} = \mathbf{I}$, $\sum_t u_{t,j}^2 = 1$, thus the equality (7) is established, moreover, since $i_1$ values from 1 to $I_1$, the corresponding $j^*$ hence goes through from 1 to $J_1$, equality (8) holds.

The Frobenius norm of $\mathcal{S}$ equals to that of $\mathcal{A}$ follows the summation from $j = 1$ to $j = J_n$ of both sides of the above equality. $\qquad\square$

# 3  The Optimization Method

In this section, we develop the optimization method for solving (2). Equation (2) can be rewritten as follows:

$$
\begin{aligned}
&\min \|\mathbf{A}_{(n)} - \mathbf{U}^{(n)}\mathbf{S}_{(n)}(\mathbf{U}^{(n+1)} \otimes \mathbf{U}^{(n+2)} \cdots \otimes \mathbf{U}^{(d)} \otimes \mathbf{U}^{(1)} \otimes \mathbf{U}^{(2)} \cdots \otimes \mathbf{U}^{(n-1)})^T\|_F^2 \\
&\text{s.t. } \mathbf{U}^{(n)} \in \mathbb{R}_+^{I_n \times J_n}, \quad \mathbf{U}^{(n)T}\mathbf{U}^{(n)} = \mathbf{I}, \quad \mathbf{S}_{(n)} \in \mathbb{R}_+^{J_n \times J_{n+1}\cdots J_d J_1 \cdots J_{n-1}}.
\end{aligned} \quad (9)
$$

where $\otimes$ denotes the Kronecker product. We remark that (9) is valid for $1 \leq n \leq d$. For simplicity, we let

$$
\mathbf{W}^{(n)} = \mathbf{S}_{(n)}(\mathbf{U}^{(n+1)} \otimes \cdots\cdots \otimes \mathbf{U}^{(d)} \otimes \mathbf{U}^{(1)} \otimes \cdots \otimes \mathbf{U}^{(n-1)})^T.
$$

Therefore, each factor matrix $\mathbf{U}^{(n)}$ can be obtained by solving the following subproblem:

$$
\begin{aligned}
&\min \|\mathbf{A}_{(n)} - \mathbf{U}^{(n)}\mathbf{W}^{(n)}\|_F^2 \\
&\text{s.t. } \mathbf{U}^{(n)} \in \mathbb{R}_+^{I_n \times J_n}, \ \mathbf{U}^{(n)T}\mathbf{U}^{(n)} = \mathbf{I}, \ \mathbf{W}^{(n)} \in \mathbb{R}_+^{J_n \times I_{n+1}\cdots I_d I_1 \cdots I_{n-1}}.
\end{aligned} \quad (10)
$$

We note that (10) is an orthogonal nonnegative matrix factorization problem. Similar to [4], we propose to solve the following optimization problem instead[1]

$$
\mathbf{U}^{(n)} = \arg\min \left\{ \|\mathbf{A}_{(n)} - \mathbf{U}^{(n)}\mathbf{U}^{(n)T}\mathbf{A}_{(n)}\|_F^2 : \mathbf{U}^{(n)} \in \mathbb{R}_+^{I_n \times J_n}, \mathbf{U}^{(n)T}\mathbf{U}^{(n)} = \mathbf{I} \right\}. \quad (11)
$$

Next we study how to solve (11) efficiently.

---

[1] Assume that $\mathcal{A}$ is orthogonally decomposable, i.e., $\mathcal{A} = \mathcal{S} \times_1 \mathbf{U}^{(1)} \cdots \times_d \mathbf{U}^{(d)}$ with $\mathbf{U}^{(n)T}\mathbf{U}^{(n)} = \mathbf{I}$ for $1 \leq n \leq d$. It is clear that $\mathbf{A}_{(n)} = \mathbf{U}^{(n)}\mathbf{W}^{(n)}$ for $1 \leq n \leq d$. Since $\mathbf{U}^{(n)}$ is orthogonal, we know that $\mathbf{W}^{(n)} = \mathbf{U}^{(n)T}\mathbf{A}_{(n)}$. It implies that $\mathbf{A}_{(n)} = \mathbf{U}^{(n)T}\mathbf{U}^{(n)}\mathbf{A}_{(n)}$. Here we propose to minimize the difference between $\mathbf{A}_{(n)}$ and $\mathbf{U}^{(n)T}\mathbf{U}^{(n)}\mathbf{A}_{(n)}$.

## 3.1 The Factor Matrix

By using Lemma 1 and Lemma 2, we rewrite problem (11) as follows: for $n = 1, 2, \ldots, d$,

$$\min_{\mathbf{K}^{(n)}} F(\mathbf{K}^{(n)}) = \frac{1}{2}\|\mathbf{A}_{(n)} - \mathbf{K}^{(n)}\mathbf{A}_{(n)}\|_F^2$$
$$\text{s.t. } tr(\mathbf{K}^{(n)}) = J_n, \quad \mathbf{K}^{(n)} = \mathbf{K}^{(n)T}, \ (\mathbf{K}^{(n)})^2 = \mathbf{K}^{(n)}, \quad \mathbf{K}^{(n)} \geq 0. \tag{12}$$

According to Lemma 2, $\mathbf{K}^{(n)}$ has a block-like structure and $\|\mathbf{K}^{(n)}\|_1 \leq I_n$. We therefore expect many entries of $\mathbf{K}^{(n)}$ are zero, and present the following convex relaxation model ,

$$\min_{\mathbf{K}^{(n)}} F(\mathbf{K}^{(n)}) = \frac{1}{2}\|\mathbf{A}_{(n)} - \mathbf{K}^{(n)}\mathbf{A}_{(n)}\|_F^2 + \theta\|\mathbf{K}^{(n)}\|_1$$
$$\text{s.t. } tr(\mathbf{K}^{(n)}) = J_n, \quad \mathbf{K}^{(n)} = \mathbf{K}^{(n)T}, \ \mathbf{0} \preceq \mathbf{K}^{(n)} \preceq \mathbf{I}, \quad \mathbf{K}^{(n)} \geq 0. \tag{13}$$

where $\|\cdot\|_1$ is $\ell_1$ norm of matrix, $\mathbf{0} \preceq \mathbf{K}^{(n)} \preceq \mathbf{I}$ denotes that matrix $\mathbf{K}^{(n)}$ and matrix $\mathbf{I} - \mathbf{K}^{(n)}$ are positive semidefinite. It is the convex hull of $\mathbf{K}^2 = \mathbf{K}$, which leads to the convex problem (13). The use of $\|\mathbf{K}^{(n)}\|_1$ is to enforce the sparsity of $\mathbf{K}^{(n)}$ and $\theta$ is a positive number to control the balance among the two terms in the objective function.

Let $\mathbf{K}^{(n)} = \mathbf{X}^{(n)}$, $\mathbf{K}^{(n)} = \mathbf{Z}^{(n)}$, $\mathbf{K}^{(n)} = \mathbf{M}^{(n)}$, we have

$$\min_{\mathbf{K}^{(n)}} F(\mathbf{K}^{(n)}) = \frac{1}{2}\|\mathbf{A}_{(n)} - \mathbf{K}^{(n)}\mathbf{A}_{(n)}\|_F^2 + \theta\|\mathbf{X}^{(n)}\|_1$$
$$\text{s.t. } \mathbf{K}^{(n)} - \mathbf{X}^{(n)} = 0, \mathbf{K}^{(n)} - \mathbf{Z}^{(n)} = 0, \mathbf{K}^{(n)} - \mathbf{M}^{(n)} = 0,$$
$$tr(\mathbf{K}^{(n)}) = J_n, \mathbf{M}^{(n)} = \mathbf{M}^{(n)T}, \mathbf{Z}^{(n)} \geq 0, \mathbf{0} \preceq \mathbf{M}^{(n)} \preceq \mathbf{I}. \tag{14}$$

We apply the alternating direction method of multipliers to solve (14). The augmented Lagrangian function of (14) is given by

$$L(\mathbf{K}^{(n)}, \mathbf{X}^{(n)}, \mathbf{Z}^{(n)}, \mathbf{M}^{(n)})$$
$$= \frac{1}{2}\|\mathbf{A}_{(n)} - \mathbf{K}^{(n)}\mathbf{A}_{(n)}\|_F^2 + \theta\|\mathbf{X}^{(n)}\|_1 + \delta_{\mathbb{R}_+^{I_n \times I_n}}(\mathbf{Z}^{(n)}) - \langle \mathbf{\Lambda}_1^{(n)}, \mathbf{K}^{(n)} - \mathbf{X}^{(n)} \rangle - \langle \mathbf{\Lambda}_2^{(n)}, \mathbf{K}^{(n)} - \mathbf{Z}^{(n)} \rangle$$
$$- \langle \mathbf{\Lambda}_3^{(n)}, \mathbf{K}^{(n)} - \mathbf{M}^{(n)} \rangle + \frac{\rho_1^{(n)}}{2}\|\mathbf{K}^{(n)} - \mathbf{X}^{(n)}\|_F^2 + \frac{\rho_2^{(n)}}{2}\|\mathbf{K}^{(n)} - \mathbf{Z}^{(n)}\|_F^2 + \frac{\rho_3^{(n)}}{2}\|\mathbf{K}^{(n)} - \mathbf{M}^{(n)}\|_F^2, \tag{15}$$

where $\delta_{\mathbb{R}_+^{I_n \times I_n}}$ denotes the indicator of $\mathbb{R}_+^{I_n \times I_n}$, i.e.,

$$\delta_{\mathbb{R}_+^{I_n \times I_n}}(\mathbf{X}) := \begin{cases} 0, & \text{if } \mathbf{X} \in \mathbb{R}_+^{I_n \times I_n}, \\ +\infty, & \text{otherwise.} \end{cases}$$

The iterative system of ADMM is given as follows:

$$(\mathbf{K}^{(n)})_{i+1} = \arg\min \left\{ L(\mathbf{K}, (\mathbf{X}^{(n)})_i, (\mathbf{Z}^{(n)})_i, (\mathbf{M}^{(n)})_i) : tr(\mathbf{K}^{(n)}) = J_n \right\},$$
$$((\mathbf{X}^{(n)})_{i+1}, (\mathbf{Z}^{(n)})_{i+1}, (\mathbf{M}^{(n)})_{i+1}) = \arg\min \left\{ L((\mathbf{K}^{(n)})_{i+1}, \mathbf{X}^{(n)}, \mathbf{Z}^{(n)}, \mathbf{M}^{(n)}) : \mathbf{0} \preceq \mathbf{M}^{(n)} \preceq \mathbf{I}, \mathbf{M}^{(n)T} = \mathbf{M}^{(n)} \right\},$$
$$(\mathbf{\Lambda}_1^{(n)})_{i+1} = (\mathbf{\Lambda}_1^{(n)})_i - \gamma^{(n)}\rho_1^{(n)}((\mathbf{K}^{(n)})_{i+1} - (\mathbf{X}^{(n)})_{i+1}),$$
$$(\mathbf{\Lambda}_2^{(n)})_{i+1} = (\mathbf{\Lambda}_2^{(n)})_i - \gamma^{(n)}\rho_2^{(n)}((\mathbf{K}^{(n)})_{i+1} - (\mathbf{Z}^{(n)})_{i+1}),$$
$$(\mathbf{\Lambda}_3^{(n)})_{i+1} = (\mathbf{\Lambda}_3^{(n)})_i - \gamma^{(n)}\rho_3^{(n)}((\mathbf{K}^{(n)})_{i+1} - (\mathbf{M}^{(n)})_{i+1}), \tag{16}$$

where $\gamma^{(n)} \in (0, (1 + \sqrt{5})/2)$.

### 3.1.1 The Computation of $\mathbf{K}^{(n)}$

For $\mathbf{K}^{(n)}$, by [8], we can solve it as

$$(\mathbf{K}^{(n)})_{i+1} = (\mathbf{B}^{(n)})_i - \frac{tr((\mathbf{B}^{(n)})_i) - J_n}{I_n}\mathbf{I}, \tag{17}$$

where

$$(\mathbf{B}^{(n)})_i = (\mathbf{P} + \mathbf{Q})\left(\mathbf{A}_{(n)}\mathbf{A}_{(n)}^T + (\rho_1^{(n)} + \rho_2^{(n)} + \rho_3^{(n)})\mathbf{I}\right)^{-1}$$

with

$$\mathbf{P} = \mathbf{A}_{(n)}\mathbf{A}_{(n)}^T + (\mathbf{\Lambda}_1^{(n)})_i + (\mathbf{\Lambda}_2^{(n)})_i + (\mathbf{\Lambda}_3^{(n)})_i, \quad \mathbf{Q} = \rho_1^{(n)}(\mathbf{X}^{(n)})_i + \rho_2^{(n)}(\mathbf{Z}^{(n)})_i + \rho_3^{(n)}(\mathbf{M}^{(n)})_i.$$

### 3.1.2 The Computation of $\mathbf{X}^{(n)}$

For $\mathbf{X}^{(n)}$, it is the shrinkage

$$(\mathbf{X}^{(n)})_{i+1} = \arg\min_{\mathbf{X}^{(n)}}\left\{\theta\|\mathbf{X}^{(n)}\|_1 + \frac{\rho_1^{(n)}}{2}\|\mathbf{X}^{(n)} - ((\mathbf{K}^{(n)})_{i+1} - \frac{1}{\rho_1^{(n)}}(\mathbf{\Lambda}_1^{(n)})_i)\|_F^2\right\}. \tag{18}$$

Thus,

$$(\mathbf{X}^{(n)})_{i+1} = \text{Shrinkage}\left((\mathbf{K}^{(n)})_{i+1} - \frac{1}{\rho_1^{(n)}}(\mathbf{\Lambda}_1^{(n)})_i, \frac{\theta}{\rho_1^{(n)}}\right), \tag{19}$$

where $\text{Shrinkage}(x, \tau) := \text{sign}(x)\max\{|x| - \tau, 0\}$ and $\text{sign}(\cdot)$ denotes the signum function, i.e.,

$$\text{sign}(x) := \left\{\begin{array}{ll} 1, & \text{if } x > 0, \\ 0, & \text{if } x = 0, \\ -1, & \text{if } x < 0. \end{array}\right.$$

### 3.1.3 The Computation of $\mathbf{Z}^{(n)}$

For $\mathbf{Z}^{(n)}$, it is the projection onto $\mathbb{R}_+^{I_n \times I_n}$,

$$(\mathbf{Z}^{(n)})_{i+1} = \arg\min_{\mathbf{Z}^{(n)}}\left\{\delta_{\mathbb{R}_+^{I_n \times I_n}}(\mathbf{Z}^{(n)}) + \frac{\rho_2^{(n)}}{2}\left\|\mathbf{Z}^{(n)} - \left((\mathbf{K}^{(n)})_{i+1} - \frac{1}{\rho_2^{(n)}}(\mathbf{\Lambda}_2^{(n)})_i\right)\right\|_F^2\right\}. \tag{20}$$

It is given by

$$(\mathbf{Z}^{(n)})_{i+1} = \Pi_{\mathbb{R}_+^{I_n \times I_n}}\left((\mathbf{K}^{(n)})_{i+1} - \frac{1}{\rho_2^{(n)}}(\mathbf{\Lambda}_2^{(n)})_i\right), \tag{21}$$

where $\Pi_{\mathbb{R}_+^{I_n \times I_n}}$ is the projection onto $\mathbb{R}_+^{I_n \times I_n}$.

### 3.1.4 The Computation of $\mathbf{M}^{(n)}$

For $\mathbf{M}^{(n)}$, it is

$$(\mathbf{M}^{(n)})_{i+1} = \arg\min\left\{\frac{\rho_3^{(n)}}{2}\left\|\mathbf{M}^{(n)} - ((\mathbf{K}^{(n)})_{i+1} - \frac{1}{\rho_3^{(n)}}(\mathbf{\Lambda}_3^{(n)})_i)\right\|_F^2 : 0 \preceq \mathbf{M}^{(n)} \preceq I, \mathbf{M}^{(n)T} = \mathbf{M}^{(n)}\right\}. \tag{22}$$

By the projection of a matrix on the symmetric positive matrix [9, Section 4.3] and [10, Lemma 2.1], we have

$$(\mathbf{M}^{(n)})_{i+1} = \frac{1}{2}\tilde{\mathbf{V}}^{(n)}\min\{\max\{\tilde{\Sigma}^{(n)},0\},1\}\tilde{\mathbf{V}}^{(n)T}, \tag{23}$$

where $\tilde{\mathbf{V}}^{(n)}$, $\tilde{\mathbf{\Sigma}}^{(n)}$ are the eigenvalue decomposition of

$$(\mathbf{K}^{(n)})_{i+1} - \frac{1}{\rho_3^{(n)}}(\mathbf{\Lambda}_3^{(n)})_i + \left((\mathbf{K}^{(n)})_{i+1} - \frac{1}{\rho_3^{(n)}}(\mathbf{\Lambda}_3^{(n)})_i\right)^T.$$

### 3.1.5 The Algorithm and Convergence Analysis

The proposed algorithm for $\{\mathbf{K}^{(n)}\}_{n=1}^d$ is given in Algorithm 1. For each $\mathbf{K}^{(n)}$ with $n = 1, 2, \cdots, d$, in the alternating direction method of multipliers, two blocks of variables $\mathbf{K}^{(n)}$ and $(\mathbf{X}^{(n)}; \mathbf{Z}^{(n)}; \mathbf{M}^{(n)})$ are updated in each iteration in our algorithm. The convergence of the algorithm can be guaranteed, see ( [8, 11, 12]). The detailed proof is given in Appendix, we follow the main idea from [8], the difference is our proof is based on our model which has one more constraint than theirs.

For $n$-th mode when $n = 1, 2, \ldots, d$,

**Theorem 1.** *Assume that $\gamma^{(n)} \in (0, (1 + \sqrt{5})/2)$, for any $\rho_1^{(n)}, \rho_2^{(n)}, \rho_3^{(n)} > 0$, the iterative sequence $((\mathbf{K}^{(n)})_i; (\mathbf{X}^{(n)})_i; (\mathbf{Z}^{(n)})_i; (\mathbf{M}^{(n)})_i; (\mathbf{\Lambda}_1^{(n)})_i; (\mathbf{\Lambda}_2^{(n)})_i; (\mathbf{\Lambda}_3^{(n)})_i)$ generated by Algorithm 1 from any initial point converges to $((\mathbf{K}^{(n)})_*; (\mathbf{X}^{(n)})_*; (\mathbf{Z}^{(n)})_*; (\mathbf{M}^{(n)})_*; (\mathbf{\Lambda}_1^{(n)})_*; (\mathbf{\Lambda}_2^{(n)})_*; (\mathbf{\Lambda}_3^{(n)})_*)$, where $((\mathbf{K}^{(n)})_*; (\mathbf{X}^{(n)})_*; (\mathbf{Z}^{(n)})_*; (\mathbf{M}^{(n)})_*)$ is a solution of (13).*

After $\mathbf{K}^{(n)}$ is computed by Algorithm 1, $\mathbf{U}^{(n)}$ can be recovered based on Lemma 2, i.e., we compute the $J_n$ eigenvectors corresponding to $J_n$ largest eigenvalues of $\mathbf{K}^{(n)}$, then use hard clustering evaluation [3] to form $\mathbf{U}^{(n)}$. Or one can simply use some known clustering methods [13–15] such as $k$-means, spectral clustering to get $\mathbf{U}^{(n)}$ from $\mathbf{K}^{(n)}$.

---

**Algorithm 1** Alternating direction method of multipliers for model ONTD

---

**Input:** Given $\mathcal{A} \in \mathbb{R}_+^{I_1 \times I_2 \times \cdots \times I_d}$, $(J_1, J_2, \cdots, J_d)$, the parameters $\theta$, $\rho_1$, $\rho_2$, $\rho_3$, $\gamma$, initial values $(\mathbf{K}^{(n)})_1 \in \mathbb{R}_+^{I_n \times I_n}$, $(\mathbf{X}^{(n)})_1 \in \mathbb{R}_+^{I_n \times I_n}$, $(\mathbf{Z}^{(n)})_1 \in \mathbb{R}_+^{I_n \times I_n}$, $(\mathbf{M}^{(n)})_1 \in \mathbb{R}_+^{I_n \times I_n}$, $(\mathbf{\Lambda}_1^{(n)})_1 \in \mathbb{R}^{I_n \times I_n}$, $(\mathbf{\Lambda}_2^{(n)})_1 \in \mathbb{R}^{I_n \times I_n}$, $(\mathbf{\Lambda}_3^{(n)})_1 \in \mathbb{R}^{I_n \times I_n}$, and the stopping criterion $\epsilon$.
**Output: K**
  1: **Step 0**. Unfolding tensor $\mathcal{A}$ from $n$ mode, obtain the $n$-mode matrix $\mathbf{A}_{(n)}$.
  2: **Step 1**. Compute $(\mathbf{K}^{(n)})_{i+1}$ by (17).
  3: **Step 2**. Compute $(\mathbf{X}^{(n)})_{i+1}$, $(\mathbf{Z}^{(n)})_{i+1}$, $(\mathbf{M}^{(n)})_{i+1}$ by (19), (21), (23), respectively.
  4: **Step 3**. Update $(\mathbf{\Lambda}_1^{(n)})_{i+1}$, $(\mathbf{\Lambda}_2^{(n)})_{i+1}$, $(\mathbf{\Lambda}_3^{(n)})_{i+1}$ by (16).
  5: **Step 4**. If the termination criterion is not met, go to Step 1.

---

## 3.2 The Core Tensor $\mathcal{S}$

After we get the factor matrices $\{\mathbf{U}^{(n)}\}_{n=1}^d$, the core tensor $\mathcal{S}$ can be obtained by

$$\mathcal{S} = \arg\min\left\{\|\mathcal{A} - \mathcal{S} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)} \cdots \times_d \mathbf{U}^{(d)}\|_F^2 : \mathcal{S} \in \mathbb{R}_+^{J_1 \times J_2 \times \cdots \times_d J_d}\right\} \tag{24}$$

Here, we use $n$-th mode $\mathbf{S}_{(n)}$ to fold into the core tensor $\mathcal{S}$. We remark that (25) is the unfolding version of (24), the solution of (25) is the same as that of (24). It is valid to choose any value of $n$ in between 1 and $d$.

It is clear that $\mathbf{S}_{(n)}$ can be solved by the following least square problem:

$$\mathbf{S}_{(n)} = \arg\min \left\{ \|\mathbf{A}_{(n)} - \mathbf{U}^{(n)}\mathbf{S}_{(n)}(\mathbf{U}^{(n+1)} \otimes \mathbf{U}^{(n+2)} \cdots \cdots \otimes \mathbf{U}^{(d)} \otimes \mathbf{U}^{(1)} \otimes \cdots \otimes \mathbf{U}^{(n-1)})^T\|_F^2 : \right.$$
$$\left. \mathbf{S}_{(n)} \in \mathbb{R}_+^{J_n \times J_{n+1} \cdots J_d J_1 \cdots J_{n-1}} \right\}. \tag{25}$$

We remark that (25) is valid for $1 \leq n \leq d$. Let $\mathbf{U}^{(\backslash n)} = \mathbf{U}^{(n+1)} \otimes \mathbf{U}^{(n+2)} \cdots \cdots \otimes \mathbf{U}^{(d)} \otimes \mathbf{U}^{(1)} \otimes \cdots \otimes \mathbf{U}^{(n-1)}$ and take use of the property of Kronecker product that $vec(\mathbf{U}^{(n)}\mathbf{S}_{(n)}\mathbf{U}^{(\backslash n)T}) = (\mathbf{U}^{(\backslash n)} \otimes \mathbf{U}^{(n)})vec(\mathbf{S}_{(n)})$, (25) can be written as follows,

$$\min_{\mathbf{S}_{(n)}} \|vec(\mathbf{A}_{(n)}) - (\mathbf{U}^{(\backslash n)} \otimes \mathbf{U}^{(n)})vec(\mathbf{S}_{(n)})\|_F^2 \quad \text{s.t.} \quad \mathbf{S}_{(n)} \in \mathbb{R}_+^{J_n \times J_{n+1} \cdots J_d J_1 \cdots J_{n-1}}. \tag{26}$$

Denote $(\mathbf{U}^{(\backslash n)} \otimes \mathbf{U}^{(n)})$ as $\mathbf{L}$, the least square solution of (26) can be obtained by

$$vec(\mathbf{S}_{(n)}) = \Pi_{\mathbb{R}_+^{I_n \times I_n}}\left((\mathbf{L}^T\mathbf{L})^{-1}\mathbf{L}^T vec(\mathbf{A}_{(n)})\right). \tag{27}$$

Core tensor $\mathcal{S}$ can be obtained by folding $\mathbf{S}_{(n)}$ from $n$-th mode.

Note that if $\{\mathbf{U}^{(n)}\}_{n=1}^d$ are column orthogonal, then

$$\mathcal{S} = \mathcal{A} \times_1 \mathbf{U}^{(1)T} \times_2 \mathbf{U}^{(2)T} \cdots \times_d \mathbf{U}^{(d)T}.$$

# 4  Experimental Results

In this section, we conduct experiments on three applied areas to test the performance of the proposed ONTD model. All experiments are run on Intel(R) Core(TM) i7-5600 CPU @2.60GHZ with 8GB of RAM.

As comparison, we use some general methods in the experiments, for example, the tensor methods such as NTD [16], HOOI [5] and matrix methods like NMF [17,18], SN-ONMF [4].

## 4.1  Feature Extraction and Face Recognition

In this subsection, we use face ORL database [19] which contains 400 images of 40 individuals. The images are captured at different times and different variations including facial details and expression. These images are in gray scale and normalized to the resolution of $112 \times 92$ pixels. We randomly select $p \times 100\%$ sample images for each person to be the training data, and the others are used for testing. Use the tensor methods, we decompose the training data with setting $(J_1, J_2)$. The slices of $\mathcal{S} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{U}^{(2)}$ can be seen as basis images. For comparison, matrix methods are also used to learn the basis matrix from the flattened training data with the same number of features. Once we get the basis matrix from the training data, the nonnegative projection of a new test sample onto each basis matrix is used as the features for recognition. The KNN classifier is used for recognition by using the extracted features, here the distance is measured by their correlations.

First we use the basis images based on 20% training data set (80 images). The number of features is setting to $J_3$. To show reconstruction capacity for original image, we define compression ratio

$$compression \ \ ratio = \frac{Compressed \ \ size}{Original \ \ size} \times 100\%$$

Table 1: The face-reconstruct average error of different methods with different number of features.

|        | $J_3 = 5$ | $J_3 = 10$ | $J_3 = 15$ | $J_3 = 20$ | $J_3 = 25$ | $J_3 = 30$ | $J_3 = 40$ | $J_3 = 60$ |
|--------|-----------|------------|------------|------------|------------|------------|------------|------------|
| ONTD   | 0.0029    | 0.0025     | 0.0023     | 0.0022     | 0.0021     | 0.0020     | 0.0019     | 0.0017     |
| HOOI   | 0.0029    | 0.0025     | 0.0023     | 0.0022     | 0.0021     | 0.0020     | 0.0019     | 0.0017     |
| NTD    | 0.0030    | 0.0027     | 0.0026     | 0.0025     | 0.0024     | 0.0024     | 0.0023     | 0.0022     |

Table 2: The face-reconstruct results of different methods by using 20% as training set.

|         | $J_3=20$ | | $J_3=40$ | | $J_3=60$ | | $J_3=80$ | |
|---------|----------|------|----------|--------|----------|--------|----------|--------|
|         | accuracy | time | accuracy | time | accuracy | time | accuracy | time |
| ONTD    | **0.7750** | 4.33 | **0.8063** | 5.93 | **0.8063** | 3.63 | 0.7531 | 5.24 |
| HOOI    | 0.7250 | 146.37 | 0.7750 | 309.97 | 0.7750 | 434.96 | 0.7063 | 350.29 |
| NTD     | 0.3000 | 121.35 | 0.3500 | 130.94 | 0.3594 | 146.22 | 0.3875 | 153.57 |
| NMF     | 0.7031 | 150.82 | 0.7156 | 176.50 | 0.7000 | 211.96 | 0.7281 | 273.69 |
| SN-ONMF | 0.7125 | 2.04 | 0.7438 | 2.06 | 0.6844 | 1.94 | **0.7688** | 1.97 |

Table 3: The face-reconstruct results of different methods by using 30% as training set.

|         | $J_3=20$ | | $J_3=40$ | | $J_3=60$ | | $J_3=80$ | |
|---------|----------|------|----------|--------|----------|--------|----------|--------|
|         | accuracy | time | accuracy | time | accuracy | time | accuracy | time |
| ONTD    | **0.8000** | 3.92 | **0.8393** | 5.10 | **0.8321** | 6.42 | **0.8250** | 6.70 |
| HOOI    | 0.7536 | 125.22 | 0.8143 | 182.70 | 0.8036 | 207.25 | 0.7321 | 387.94 |
| NTD     | 0.4750 | 149.29 | 0.4571 | 157.66 | 0.4571 | 163.61 | 0.4321 | 214.32 |
| NMF     | 0.7357 | 219.93 | 0.7893 | 256.74 | 0.7536 | 291.91 | 0.7321 | 372.42 |
| SN-ONMF | 0.7321 | 3.99 | 0.7857 | 4.06 | 0.7321 | 4.11 | 0.6500 | 4.24 |

and the average reconstruction error

$$avg.error = \frac{1}{N} \sum_{k=1}^{N} \frac{\|I_{original}^k - I_{recon}^k\|_F}{\|I_{original}^k\|_F},$$

where $I_{original}^k$ represents the $k$-th original image, $I_{recon}^k$ represents the $k$-th reconstruction image, $N$ is the number of the images.

In Table 1, for tensor methods, we summarize the average reconstruction error according to different $(J_1, J_2, J_3)$, here $(J_1, J_2) = (5, 5)$, $(10, 10)$, $(15, 15)$, $(20, 20)$, $(25, 25)$, $(30, 30)$, $(40, 40)$, $(60, 60)$. It can be seen from Table 1 that the reconstruction error decrease when we increase $J_3$. From Table 1, the reconstruction error of ONTD and HOOI are similar, which is less than NTD.

In the following, we set $(J_1, J_2)$ to be $(15, 15)$ and get the basic matrix from different training data set. Test the the recognition accuracy on the rest data set (i.e., testing set) according to the basic matrix. We show the recognition results on Table 2-Table 5. It can be seen that the number of features ($J_3$) affects the recognition. A Large value of $J_3$ often leads to a higher accuracy, but when $J_3$ is too large, the accuracy decrease. We also notice that the recognition accuracy can be increase when using more training data. From Table 2-Table 5, the accuracy obtained by our method (ONTD) is higher than the other methods in most cases and ONTD takes much less time than the other tensor methods.

Table 4: The face-reconstruct results of different methods via using 40% as training set.

|  | $J_3=20$ | | $J_3=40$ | | $J_3=60$ | | $J_3=80$ | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
|  | accuracy | time | accuracy | time | accuracy | time | accuracy | time |
| ONTD | **0.8750** | 10.06 | **0.8750** | 9.12 | **0.8875** | 9.84 | **0.8958** | 8.46 |
| HOOI | 0.8083 | 186.10 | 0.8375 | 309.56 | 0.8500 | 416.10 | 0.8500 | 542.87 |
| NTD | 0.4417 | 215.10 | 0.4667 | 220.56 | 0.4833 | 267.85 | 0.5542 | 303.57 |
| NMF | 0.7667 | 376.67 | 0.8333 | 420.3 | 0.8250 | 457.10 | 0.7833 | 502.94 |
| SN-ONMF | 0.8208 | 6.53 | 0.8208 | 7.42 | 0.7917 | 6.83 | 0.7500 | 6.46 |

Table 5: The face-reconstruct results of different methods by using 50% as training set.

|  | $J_3=20$ | | $J_3=40$ | | $J_3=60$ | | $J_3=80$ | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
|  | accuracy | time | accuracy | time | accuracy | time | accuracy | time |
| ONTD | **0.8800** | 12.24 | **0.8800** | 12.91 | **0.8900** | 5.74 | **0.8800** | 11.60 |
| HOOI | 0.8200 | 186.47 | 0.8400 | 296.72 | 0.8600 | 412.33 | 0.8100 | 552.86 |
| NTD | 0.4450 | 296.75 | 0.5250 | 328.19 | 0.4300 | 347.22 | 0.4950 | 375.75 |
| NMF | 0.8050 | 493.78 | 0.8150 | 470.28 | 0.7950 | 586.04 | 0.8050 | 684.95 |
| SN-ONMF | 0.8350 | 9.68 | 0.8400 | 9.81 | 0.8250 | 9.92 | 0.7600 | 10.17 |

## 4.2  Image Representation

In the following data sets, the image sets are belong to several different classes, each class is represented by a 3-order tensor. For example, there are $N$ tensor objects size of $I_1 \times I_2 \times I_3$ that can be classified into $r$ classes. To implement the tensor methods, we concatenate these $N$ 3-order tensor objects to a $I_1 \times I_2 \times I_3 \times N$ tensor. We follow the idea in [20, 21], i.e., given approximate rank $(J_1, J_2, J_3)$, use the tensor methods to find three common projection matrices $\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}$ for these 3-order tensor objects. These tensor objects then can be represented as $N$ corresponding core tensors of size $(J_1, J_2, J_3)$. When use the matrix methods, we vectorize $N$ tensor objects and get a $I_1 I_2 I_3 \times N$ matrix first.

To test the performance of the proposed model, we reduce the dimension of tensor objects and report the space savings.

$$Space \ \ savings = 1 - \frac{Compressed \ \ size}{Original \ \ size}$$

Since the data sets have ground truth class label, We also classify these $N$ reduced tensor objects (core tensor) with nearest neighborhood classifier and use the leave-one-out scheme, then show the average precision to evaluate the classification performance.

### 4.2.1  The ORL Database of Faces

Here we use ORL database of faces again. As we know, there are ten different images of each of 40 distinct individual. Each image is of $92 \times 112$ pixels. We first uniformly resized them into $20 \times 20$, then randomly choose two images from ten for each subjects to form a $20 \times 20 \times 2$ tensor. Therefore, for each person, we get five tensor objects. From 40 distinct individuals, we can obtain 200 tensor objects totally. We select 40% tensor objects from 200 tensor objects randomly to form a $20 \times 20 \times 2 \times 80$ training tensor to tune a good approximate rank. The rest 120 tensor objects are used to form the test data set to evaluate the results.

Firstly we set the approximate rank to be $[J_1, J_2, J_3]$ for the tensor methods, and apply these methods on the training set, we get the common matrices $(\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)})$ and 80 core tensors

of $[J_1, J_2, J_3]$. For matrix methods, we unfold the tensors to a $800 \times 80$ matrix. The approximate rank for matrix is set to be $r$. We compute the best parameters $(J_1, J_2, J_3)$ to get the highest precision by using the the tensor methods, or the best parameter $r$ to get its largest precision by each matrix method. Then the parameters $(J_1, J_2, J_3)$ or $r$ are used in the test set. We show the results in Table 6.

From Table 6, compare to other four methods, ONTD has the highest precision and the least time.

Table 6: The clustering results for ORL database.

|         | Best parameter | precision | Space saving | time  |
|---------|----------------|-----------|--------------|-------|
| ONTD    | $[9, 3, 2]$    | 0.1375    | 0.9300       | 0.03  |
| HOOI    | $[3, 3, 2]$    | 0.125     | 0.9762       | 18.79 |
| NTD     | $[7, 3, 2]$    | 0.0917    | 0.9454       | 5.76  |
| NMF     | 15             | 0.125     | 0.8562       | 0.15  |
| SN-ONMF | 25             | 0.1       | 0.7604       | 6.366 |

### 4.2.2  Cambridge Gesture database

Cambridge Gesture data sets contains 900 images from nine hand gesture classes [22]. Each class has 100 image sequences (5 illuminations $\times$ 10 arbitrary motions $\times$ 2 subjects). Each image has $320 \times 240$ pixels. For all video sequences, we use the gray scale representation and uniformly resized them into $20 \times 20 \times 32$. We get 900 tensor objects of size $20 \times 20 \times 32$. We select 30% tensor objects from 900 tensor objects randomly, i.e., 270 tensor objects of size $20 \times 20 \times 32$. From 270 tensor objects, we randomly use 90 tensor objects to form a $20 \times 20 \times 32 \times 90$ training tensor to tune a good approximate rank for all methods, a $20 \times 20 \times 32 \times 180$ tensor which is formed by the rest 180 tensor objects, are utilized as the test set to evaluate the results.

We set the approximate rank to be $(J_1, J_2, J_3)$ for all the tensor methods and apply them on the training set, we get the common matrices $(\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)})$ and 90 core tensors of $(J_1, J_2, J_3)$. For matrix methods, similarly, we unfold the tensors to a $12800 \times 90$ matrix. The approximate rank for matrix is set to be $r$. We tune the best parameters, i.e., approximate rank $(J_1, J_2, J_3)$ or rank $r$ by training set. The best parameters are then used in the test set.

The results of test set are listed in Table 7. From space savings, we find that matrix methods needs more storage than tensor methods, it implies that tensor methods performs better in the data reduction. It is very important especially when data set is extreme large. Note that the precision of ONTD is the largest and the time is the least of all methods.

Table 7: The clustering results for Cambridge database.

|         | Best parameter | precision | Space saving | time  |
|---------|----------------|-----------|--------------|-------|
| ONTD    | $[7, 7, 3]$    | 0.8889    | 0.9884       | 0.04  |
| HOOI    | $[7, 7, 2]$    | 0.8667    | 0.9922       | 13.64 |
| NTD     | $[8, 7, 3]$    | 0.6333    | 0.9867       | 93.14 |
| NMF     | 10             | 0.7556    | 0.9437       | 3.97  |
| SN-ONMF | 25             | 0.6444    | 0.8592       | 12.64 |

Table 8: The clustering results for KTH database.

|         | Best parameter | precision | Space saving | time  |
|---------|----------------|-----------|--------------|-------|
| ONTD    | [3, 6, 2]      | 0.4400    | 0.9970       | 0.02  |
| HOOI    | [2, 2, 2]      | 0.3500    | 0.9993       | 3.94  |
| NTD     | [2, 3, 5]      | 0.3000    | 0.9975       | 51.52 |
| NMF     | 3              | 0.3800    | 0.9698       | 0.53  |
| SN-ONMF | 16             | 0.4200    | 0.8388       | 2.80  |

### 4.2.3 KTH Human Action database

KTH database [23] contains six types of human actions, including walking, running, jogging, boxing, hand clapping and hand waving. There are totally 600 videos and each type has 100 videos. The action videos are performed several times by 25 subjects in 4 different scenario (outdoors, outdoors with scale variation, outdoors with different clothes and indoors). The resolution of video frames is $160 \times 120$ pixels. The frames for the first time action of each video are extracted as our video data. To standardize the length of videos, we sample 64 frames from each video and take the middle 32 frames. For each video, we use grayscale representation and resize into $20 \times 20 \times 32$. We get 600 tensor objects whose size is $20 \times 20 \times 32$, then we select 200 tensor objects randomly. Of 200 selected tensor objects, 100 is used to form a $20 \times 20 \times 32 \times 100$ training tensor, the rest 100 tensor objects are utilized as the test set to show the performance of these methods.

For each tensor object, we set the size of the core tensor to be $[J_1, J_2, J_3]$, and apply the methods on the training set. For matrix methods, similarly, we unfold the tensor to be a $12800 \times 100$ matrix and the approximate rank is set to be $r$. We tune the best approximate rank for each method that computed on the training set to get the largest precision. Then use these parameters on test set to test the performance of each method. From Table 8, the precision of ONTD is the largest of all methods and the time cost is the least too. From space saving, it has a very similar results for each tensor method. All tensor methods perform much better than matrix methods in data reduction.

### 4.2.4 The MNIST database

The MNIST database of handwritten digits has a training set of 60000 examples, and a test set of 10000 examples. There are ten different patterns from 0 to 9. Each image is $28 \times 28$ pixels. Of the training data set, we randomly choose 90% images of each pattern, i.e. 5400 images, to form 54 tensor objects whose size is $28 \times 28 \times 100$ tensor. We obtain 540 tensor objects totally from 10 different patterns. Therefore the size of the training tensor is $28 \times 28 \times 100 \times 540$. For test data, similarly, from test set, we select 80% images of each pattern, i.e. 800 images, to form 8 tensor objects whose size is $28 \times 28 \times 100$. Thus, there are totally 80 tensor objects from 10 patterns. We get the testing tensor sizes of $28 \times 28 \times 100 \times 80$ finally.

We set the approximate rank to be $[J_1, J_2, J_3]$ for all the tensor methods and $r$ for all the matrix methods, then apply these methods on the training data. We tune best parameters approximate rank $(J_1, J_2, J_3)$ or rank $r$. Then the best parameter are applied in the test data. In Table 9, we list the results of test data. It is clearly that all tensor methods perform much better than matrix methods in data reduction. The precision of the methods are all very high, and ONTD takes much less time than other tensor methods.

## 4.3 Hyperspectral Unmixing

The hyperspectral imaging collects information from the object by taking at different wavelengths. The images are obtained by measuring the percentage of the light hitting a material which is called

Table 9: The clustering results for MNIST database.

| | Best parameter | precision | Space saving | time |
|---|---|---|---|---|
| ONTD | $[2, 2, 4]$ | 1 | 0.9997 | 0.31 |
| HOOI | $[2, 2, 3]$ | 1 | 0.9998 | 103.68 |
| NTD | $[2, 2, 4]$ | 1 | 0.9997 | 656.86 |
| NMF | 5 | 0.9875 | 0.9374 | 17.17 |
| SN-ONMF | 6 | 1 | 0.9249 | 0.18 |

reflectance. Like other spectral imaging, the purpose of hyperspectral imaging is to find objects, identify materials or detect processes. It therefore has wide applications in agriculture, mineralogy, physics, environment and many other fields. Hyperspectral unmixing aims to classify the pixels to different clusters, with each corresponding to a material.

We apply the partial ONTD model on hyperspectral unmixing, precisely, given an $I_1 \times I_2 \times I_3$ nonnegative tensor $\mathcal{A}$ and a factorization rank $r$, solve

$$\min_{\mathbf{U}^{(1)}, \mathcal{S}} \|\mathcal{A} - \mathcal{S} \times_1 \mathbf{U}^{(1)} \times_2 \mathbf{I} \times_3 \mathbf{I}\|_F^2$$

$$\text{s.t. } \mathcal{S} \geq 0, \ \mathbf{U}^{(1)} \geq 0, \ \mathbf{U}^{(1)T}\mathbf{U}^{(1)} = \mathbf{I},$$

$$\mathbf{U}^{(1)} \in \mathbb{R}^{I_1 \times r}, \mathcal{S} \in \mathbb{R}^{r \times I_2 \times I_3}.$$

In this subsection, we use the Samson data set [24]. In the image, there are $952 \times 952$ pixels, and each pixel is recorded at 156 channels that cover the wavelengths from 401 nm to 889 nm. We use a region of $95 \times 95$ pixels starting from $(252, 332)$ pixel in the original image. Three different materials are in this image, they are "Tree", "Rock", and "Water" respectively. We form a 3-order tensor $\mathcal{A} \in \mathbb{R}^{156 \times 95 \times 95}$, where 156 represents the number of spectral bands, 95 and 95 denote the row and column number of the hypercube, respectively. Moreover, we set $r = 3$ here.

In this example, the aim of NTD is finding three factor matrices and one core tensor. We do a minor revision on NTD, i.e., using their method to find the 1st factor matrix size of $156 \times 3$ and the core tensor size of $3 \times 95 \times 95$, we refer it to NTD1 here. Now we use tensor methods (ONTD, NTD, NTD1) on tensor $\mathcal{A}$, then get a $156 \times 3$ factor matrix $\mathbf{U}^{(1)}$ and a $3 \times 95 \times 95$ factor tensor $\mathcal{S}$. The $i$-th feature is obtained by hard clustering [3] base on the first array of tensor $\mathcal{S}$, where $i = 1, 2, 3$. When we use matrix methods (NMF, SN-ONMF), tensor $\mathcal{A}$ should be firstly reshaped to a pixels $\times$ spectral matrix $\mathbf{A}$ whose size is $9025 \times 156$. The factor matrices size of $9025 \times 3$ and $3 \times 156$ can be obtained by the matrix methods. After hard clustering is used on the $9025 \times 3$ matrix, we can obtain the $i$-th feature image by reshaping its $i$-th column to a $95 \times 95$ matrix. Worth to say, different from [4], here we impose the orthogonal constraint on spectral-class matrix whose size is $156 \times 3$. Note that we don't apply HOOI because the nonnegative value makes it difficult shown in unmixing.

The groundtruth and the numerical results are displayed in Figure 1. According to the results, our method displays a good clustering performance. The tree, rock and water are extracted. But for the other methods, they do not perform well because they are not able to separate these materials well.

Table 10: The hyperspectral unmixing results of different methods.

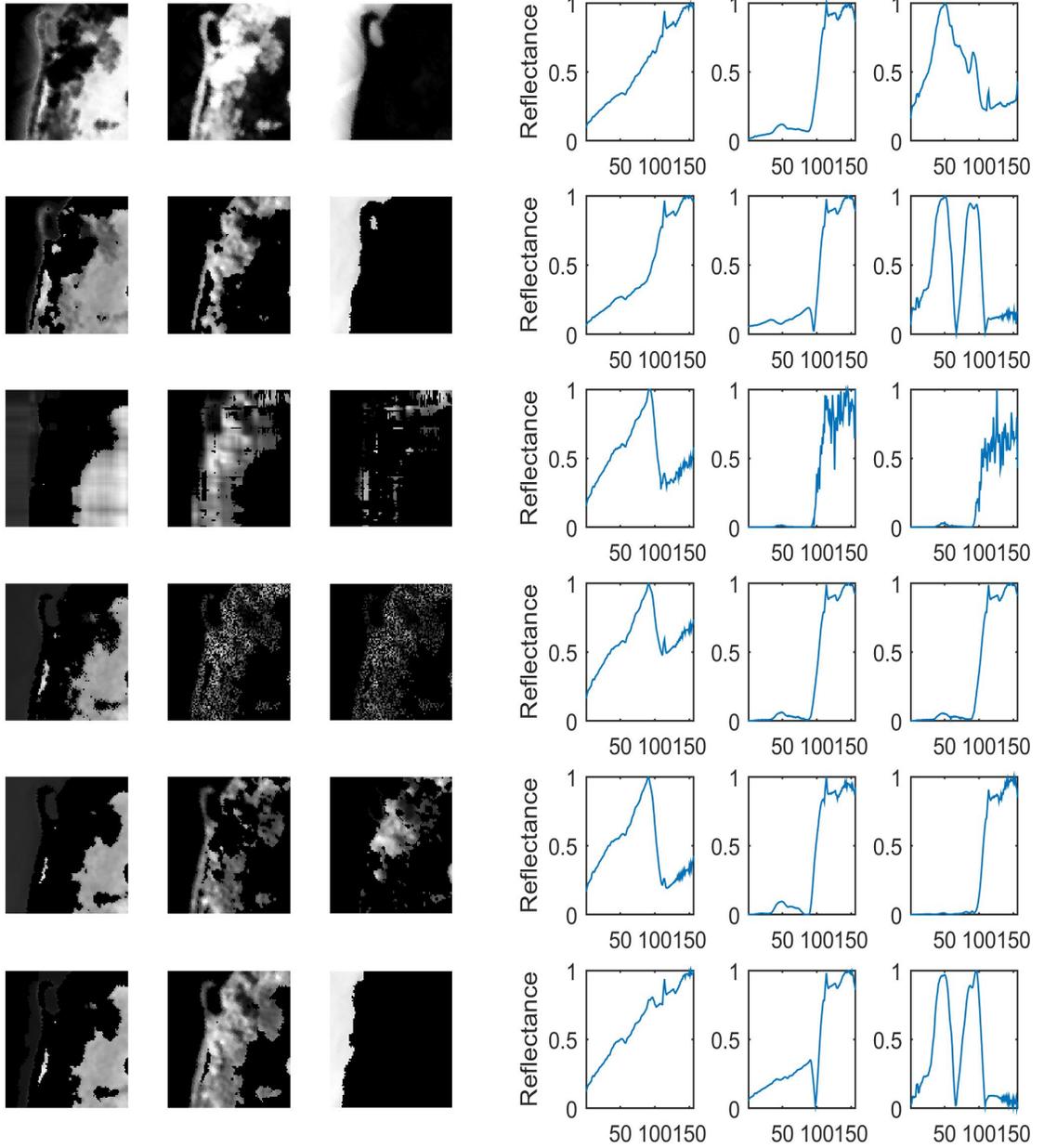| | ONTD | NTD | NTD1 | NMF | SN-ONMF |
|---|---|---|---|---|---|
| similarity | **0.9083** | 0.5674 | 0.5237 | 0.5142 | 0.8960 |
| comp. time | **2.62**s | 82.03s | 16.88s | 13.87s | 5.47s |

Figure 1: Left: Rock, Tree, Water; Right: Reflectance of Rock, Tree, Water. From the top to bottom: groundtruth, ONTD, NTD, NTD1, NMF, SN-ONMF.

To evaluate the result, in this case, we define the following similarity metric:

$$similarity = \frac{1}{r} \sum_{i=1}^{r} \frac{< \mathbf{h}_i, \mathbf{g}_i >}{\|\mathbf{h}_i\|_2 \|\mathbf{g}_i\|_2}$$

where $\{\mathbf{h}_i\}_{i=1}^{r}$ are the extracted features, $\{\mathbf{g}_i\}_{i=1}^{r}$ are the groundtruth features. It describes the

15

similarity of the groundtruth feature space and the computed feature space. Larger value implies a better result. The clustering similarity and computational time (in seconds) of the above results are list in Table 10. Compare to other methods, our model shows the highest similarity and it is the fastest. One can also find that NTD takes more time than the other methods. For matrix method, although SNONMF takes 5.47s, the similarity is high at 0.8960, only less than ONTD.

## 5 Conclusion

In this paper we have studied the orthogonal nonnegative Tucker decomposition problem and developed a structured convex optimization algorithm. The convergence analysis is given. We employ ONTD on the image data sets from the real world applications including face recognition, image representation, hyperspectral unmixing. It shows a good performance.

## Appendix

In the following, we give the proof of Theorem 1. For simplicity, from the algorithm, we give a general model as follows,

$$\min F(X, Y, Z, K) = \frac{1}{2}\|A - KA\|_F^2 + \theta\|X\|_1$$

$$\text{s.t.} \quad \begin{pmatrix} I \\ I \\ I \end{pmatrix} K + \begin{pmatrix} -I & 0 & 0 \\ 0 & -I & 0 \\ 0 & 0 & -I \end{pmatrix} \begin{pmatrix} X \\ Z \\ M \end{pmatrix} = 0,$$

$$tr(K) = J, \ M^T = M, \ 0 \preceq M \preceq I, \ Z \geq 0.$$

It can be written as

$$\min F(K, M, X, Z) = f(K) + h(M) + g(X) + \delta(Z)$$

$$\text{s.t.} \quad \begin{pmatrix} I \\ I \\ I \end{pmatrix} K + \begin{pmatrix} -I & 0 & 0 \\ 0 & -I & 0 \\ 0 & 0 & -I \end{pmatrix} \begin{pmatrix} X \\ Z \\ M \end{pmatrix} = 0,$$

where $f(K) = \frac{1}{2}\|A - KA\|_F^2 + f_1(K)$, and $f_1(K) = \begin{cases} 0, & \text{if } tr(K) = J; \\ +\infty, & \text{otherwise.} \end{cases}$, $g(X) = \theta\|X\|_1$, $h(M) = \begin{cases} 0, & \text{if } M^T = M, \ 0 \preceq M \preceq I; \\ +\infty. & \text{otherwise.} \end{cases}$, $\delta(Z) = \begin{cases} 0, & \text{if } Z \in \mathbb{R}_+; \\ +\infty, & \text{otherwise.} \end{cases}$. The general model can be simply written as

$$\min F = f(K) + h(M) + g(X) + \delta(Z) \quad \text{s.t.} \quad K = X, \ K = Z, \ K = M. \tag{28}$$

The algorithm for the general model is in the following:

$$K^i = \arg\min\{f(K) + \frac{\lambda_1}{2}\|K - M^{i-1} + b_1^{i-1}\|_F^2 + \frac{\lambda_2}{2}\|K - X^{i-1} + b_2^{i-1}\|_F^2 + \frac{\lambda_3}{2}\|K - Z^{i-1} + b_3^{i-1}\|_F^2\};$$

$$M^i = \arg\min\{h(M) + \frac{\lambda_1}{2}\|K^i - M + b_1^{i-1}\|_F^2\};$$

$$X^i = \arg\min\{g(X) + \frac{\lambda_2}{2}\|K^i - X + b_2^{i-1}\|_F^2\};$$

$$Z^i = \arg\min\{\delta(Z) + \frac{\lambda_3}{2}\|K^i - Z + b_3^{i-1}\|_F^2\};$$

$$b_1^i = b_1^{i-1} + K^i - M^i;$$

$$b_2^i = b_2^{i-1} + K^i - X^i;$$

$$b_3^i = b_3^{i-1} + K^i - Z^i.$$

The lagrangian can be written as

$$
\begin{aligned}
L =& f(K) + h(M) + g(X) + \delta(Z) + \frac{\lambda_1}{2}\|K - M\|_F^2 + \lambda_1\langle K - M, b_1\rangle + \frac{\lambda_2}{2}\|K - X\|_F^2 \\
&+ \lambda_2\langle K - X, b_2\rangle + \frac{\lambda_3}{2}\|K - Z\|_F^2 + \lambda_3\langle K - Z, b_3\rangle.
\end{aligned}
\tag{29}
$$

**Definition 1.** $(K^*, M^*, X^*, Z^*, b_1^*, b_2^*, b_3^*)$ *is a saddle point if*

$$L(K^*, M^*, X^*, Z^*, b_1, b_2, b_3) \le L(K^*, M^*, X^*, Z^*, b_1^*, b_2^*, b_3^*) \le L(K, M, X, Z, b_1^*, b_2^*, b_3^*) \tag{30}$$

*for any* $(K, M, X, Z, b_1, b_2, b_3)$.

**Lemma 3.** $(K^*, M^*, X^*, Z^*)$ *is a solution of problem (28) if and only if there exist* $b_1^*$, $b_2^*$, $b_3^*$ *such that* $(K^*, M^*, X^*, Z^*, b_1^*, b_2^*, b_3^*)$ *is a saddle point.*

*Proof.* $(K^*, M^*, X^*, Z^*, b_1^*, b_2^*, b_3^*)$ is a saddle point. From

$$L(K^*, M^*, X^*, Z^*, b_1, b_2, b_3) \le L(K^*, M^*, X^*, Z^*, b_1^*, b_2^*, b_3^*),$$

we get, $\forall b_1, b_2, b_3$,

$$\lambda_1\langle K^* - M^*, b_1\rangle + \lambda_2\langle K^* - X^*, b_2\rangle + \lambda_3\langle K^* - Z^*, b_3\rangle \le \lambda_1\langle K^* - M^*, b_1^*\rangle + \lambda_2\langle K^* - X^*, b_2^*\rangle + \lambda_3\langle K^* - Z^*, b_3^*\rangle.$$

Let $b_1 = b_1^*$, $b_2 = b_2^*$, $b_3 = b_3^* \pm \triangle b_3$, we get $K^* = Z^*$. Similarly, $K^* = M^*$, $K^* = X^*$. Also $\forall (K, M, X, Z)$,

$$L(K^*, M^*, X^*, Z^*, b_1^*, b_2^*, b_3^*) \le L(K, M, X, Z, b_1^*, b_2^*, b_3^*),$$

let $K = M = X = Z$, we have

$$f(K^*) + h(M^*) + g(X^*) + \delta(Z^*) \le f(K) + h(M) + g(X) + \delta(Z),$$

i.e., $(K^*, M^*, X^*, Z^*)$ is minimizer of $F(K, M, X, Z)$.

If $(K^*, M^*, X^*, Z^*)$ is solution of (28), i.e., $K^* = M^* = X^* = Z^*$, then the left inequality of (30) established. Moreover, $\exists b_1^*, b_2^*, b_3^*$ such that

$$-\lambda_1 b_1^* - \lambda_2 b_2^* - \lambda_3 b_3^* \in \partial f(K^*), \quad \lambda_1 b_1^* \in \partial h(M^*), \quad \lambda_2 b_2^* \in \partial g(X^*), \quad \lambda_3 b_3^* \in \partial\delta(Z^*).$$

For $\forall M, X, Z$,

$$f(K) - f(K^*) \ge \nabla f(K^*)(K - K^*), \quad h(M) - h(M^*) \ge \nabla h(M^*)(M - M^*),$$

$$g(X) - g(X^*) \geq \nabla g(X^*)(X - X^*), \quad \delta(Z) - \delta(Z^*) \geq \nabla \delta(Z^*)(Z - Z^*),$$

i.e.,

$$
\begin{aligned}
f(K^*) &\leq f(K) + \langle \lambda_1 b_1^* + \lambda_2 b_2^* + \lambda_3 b_3^*, K - K^* \rangle, \\
h(M^*) &\leq h(M) - \langle \lambda_1 b_1^*, M - M^* \rangle, \\
g(X^*) &\leq g(X) - \langle \lambda_2 b_2^*, X - X^* \rangle, \\
\delta(Z^*) &\leq \delta(Z) - \langle \lambda_3 b_3^*, Z - Z^* \rangle,
\end{aligned}
\tag{31}
$$

the summation of above inequality (31) yields the right inequality of (30). $\qquad \square$

**Theorem 2.** $\{(K^i, M^i, X^i, Z^i)\}$ *generated by algorithm from any starting point converges to a minimum of (28).*

*Proof.* Let $(K^*, M^*, X^*, Z^*)$ be an optimal solution, from algorithm,

$$
\begin{aligned}
K^* &= \arg\min L(K, M^*, X^*, Z^*, b_1^*, b_2^*, b_3^*); \\
M^* &= \arg\min L(K^*, M, X^*, Z^*, b_1^*, b_2^*, b_3^*); \\
X^* &= \arg\min L(K^*, M^*, X, Z^*, b_1^*, b_2^*, b_3^*); \\
Z^* &= \arg\min L(K^*, M^*, X^*, Z, b_1^*, b_2^*, b_3^*);
\end{aligned}
$$

for any $K, M, X, Z \in \mathbb{R}^{m \times m}$. From (31), we have,

$$
\begin{aligned}
f(K) - f(K^*) + \lambda_1 \langle K - K^*, b_1^* \rangle + \lambda_2 \langle K - K^*, b_2^* \rangle + \lambda_3 \langle K - K^*, b_3^* \rangle \\
= f(K) - f(K^*) + \lambda_1 \langle K - K^*, K^* - M^* + b_1^* \rangle \\
+ \lambda_2 \langle K - K^*, K^* - X^* + b_2^* \rangle + \lambda_3 \langle K - K^*, K^* - Z^* + b_3^* \rangle \geq 0, \\
h(M) - h(M^*) + \lambda_1 \langle M^* - K^* - b_1^*, M - M^* \rangle \geq 0, \\
g(X) - g(X^*) + \lambda_2 \langle X^* - K^* - b_2^*, X - X^* \rangle \geq 0, \\
\delta(Z) - \delta(Z^*) + \lambda_3 \langle Z^* - K^* - b_3^*, Z - Z^* \rangle \geq 0.
\end{aligned}
\tag{32}
$$

Because of the algorithm, $(K^i, M^i, X^i, Z^i)$ for any $K, M, X, Z$, we have

$$f(K) - f(K^i) \geq \nabla f^T(K^i)(K - K^i),$$

where

$$\nabla f^T(K^i) = -\lambda_1(K^i - M^{i-1} + b_1^{i-1}) - \lambda_2(K^i - X^{i-1} + b_2^{i-1}) - \lambda_3(K^i - Z^{i-1} + b_3^{i-1}).$$

So,

$$
\begin{aligned}
f(K) &- f(K^i) + \lambda_1 \langle K^i - M^{i-1} + b_1^{i-1}, K - K^i \rangle \\
&+ \lambda_2 \langle K^i - X^{i-1} + b_2^{i-1}, K - K^i \rangle + \lambda_3 \langle K^i - Z^{i-1} + b_3^{i-1}, K - K^i \rangle \geq 0.
\end{aligned}
\tag{33}
$$

Similar, we get

$$
\begin{aligned}
h(M) - h(M^i) + \lambda_1 \langle M^i - K^i - b_1^{i-1}, M - M^i \rangle &\geq 0, \\
g(X) - g(X^i) + \lambda_2 \langle X^i - K^i - b_2^{i-1}, X - X^i \rangle &\geq 0, \\
\delta(Z) - \delta(Z^i) + \lambda_3 \langle Z^i - K^i - b_3^{i-1}, Z - Z^i \rangle &\geq 0.
\end{aligned}
\tag{34}
$$

Let $K = K^i$ in (32) and $K = K^*$ in (33), then

$$f(K^i)-f(K^*)+\lambda_1\langle K^i-K^*, K^*-M^*+b_1^*\rangle+\lambda_2\langle K^i-K^*, K^*-X^*+b_2^*\rangle+\lambda_3\langle K^i-K^*, K^*-Z^*+b_3^*\rangle \geq 0,$$

$$f(K^*) - f(K^i) + \lambda_1\langle K^* - K^i, K^i - M^{i-1} + b_1^{i-1}\rangle + \lambda_2\langle K^* - K^i, K^i - X^{i-1} + b_2^{i-1}\rangle$$
$$+ \lambda_3\langle K^* - K^i, K^i - Z^{i-1} + b_3^{i-1}\rangle \geq 0.$$

The summation of above inequality is

$$\lambda_1\langle K^i - K^*, K^* - M^* + b_1^* - K^i + M^{i-1} - b_1^{i-1}\rangle + \lambda_2\langle K^i - K^*, K^* - X^* + b_2^* - K^i + X^{i-1} - b_2^{i-1}\rangle$$
$$+\lambda_3\langle K^i - K^*, K^* - Z^* + b_3^* - K^i + Z^{i-1} - b_3^{i-1}\rangle \geq 0.$$

Similarly, we get

$$\lambda_1\langle M^i - M^*, M^* - K^* - b_1^* - M^i + K^i + b_1^{i-1}\rangle \geq 0,$$
$$\lambda_2\langle X^i - X^*, X^* - K^* - b_2^* - X^i + K^i + b_2^{i-1}\rangle \geq 0,$$
$$\lambda_3\langle Z^i - Z^*, Z^* - K^* - b_3^* - Z^i + K^i + b_3^{i-1}\rangle \geq 0,$$

Let $\triangle K^i = K^i - K^*$, $\triangle M^i = M^i - M^*$, $\triangle X^i = X^i - X^*$, $\triangle Z^i = Z^i - Z^*$, $\triangle b_1^i = b_1^i - b_1^*$, $\triangle b_2^i = b_2^i - b_2^*$, $\triangle b_3^i = b_3^i - b_3^*$, so

$$\lambda_1\langle \triangle K^i, \triangle M^{i-1} - \triangle K^i - \triangle b_1^{i-1}\rangle + \lambda_2\langle \triangle K^i, \triangle X^{i-1} - \triangle K^i - \triangle b_2^{i-1}\rangle$$
$$+\lambda_3\langle \triangle K^i, \triangle Z^{i-1} - \triangle K^i - \triangle b_3^{i-1}\rangle \geq 0,$$
$$\lambda_1\langle \triangle M^i, \triangle K^i - \triangle M^i + \triangle b_1^{i-1}\rangle \geq 0,$$
$$\lambda_2\langle \triangle X^i, \triangle K^i - \triangle X^i + \triangle b_2^{i-1}\rangle \geq 0, \qquad (35)$$
$$\lambda_3\langle \triangle Z^i, \triangle K^i - \triangle Z^i + \triangle b_3^{i-1}\rangle \geq 0.$$

Summation of the above inequalities, we get that

$$\lambda_1(\langle \triangle M^i - \triangle K^i, \triangle b_1^{i-1}\rangle + \langle \triangle M^i, \triangle K^i - \triangle M^i\rangle + \langle \triangle K^i, \triangle M^{i-1} - \triangle K^i\rangle)$$
$$+\lambda_2(\langle \triangle X^i - \triangle K^i, \triangle b_2^{i-1}\rangle + \langle \triangle X^i, \triangle K^i - \triangle X^i\rangle + \langle \triangle K^i, \triangle X^{i-1} - \triangle K^i\rangle)$$
$$+\lambda_3(\langle \triangle Z^i - \triangle K^i, \triangle b_3^{i-1}\rangle + \langle \triangle Z^i, \triangle K^i - \triangle Z^i\rangle + \langle \triangle K^i, \triangle Z^{i-1} - \triangle K^i\rangle) \geq 0,$$

i.e.,

$$\lambda_1\langle \triangle M^i - \triangle K^i, \triangle b_1^{i-1}\rangle - \lambda_1\|\triangle K^i - \triangle M^i\|^2 - \lambda_1\langle \triangle M^i - \triangle M^{i-1}, \triangle K^i\rangle$$
$$+\lambda_2\langle \triangle X^i - \triangle K^i, \triangle b_2^{i-1}\rangle - \lambda_2\|\triangle K^i - \triangle X^i\|^2 - \lambda_2\langle \triangle X^i - \triangle X^{i-1}, \triangle K^i\rangle$$
$$+\lambda_3\langle \triangle Z^i - \triangle K^i, \triangle b_3^{i-1}\rangle - \lambda_3\|\triangle K^i - \triangle Z^i\|^2 - \lambda_3\langle \triangle Z^i - \triangle Z^{i-1}, \triangle K^i\rangle \geq 0.$$

On the other hand,

$$\begin{aligned} \triangle b_1^i &= \triangle b_1^{i-1} + \triangle K^i - \triangle M^i, \\ \triangle b_2^i &= \triangle b_2^{i-1} + \triangle K^i - \triangle X^i, \\ \triangle b_3^i &= \triangle b_3^{i-1} + \triangle K^i - \triangle Z^i. \end{aligned}$$

19

We get that,

$$\lambda_1(\|\triangle b_1^{i-1}\|^2 - \|\triangle b_1^i\|^2) + \lambda_2(\|\triangle b_2^{i-1}\|^2 - \|\triangle b_2^i\|^2) + \lambda_3(\|\triangle b_3^{i-1}\|^2 - \|\triangle b_3^i\|^2) \tag{36}$$
$$= \lambda_1(-2\langle \triangle K^i - \triangle M^i, \triangle b_1^{i-1}\rangle - \|\triangle K^i - \triangle M^i\|^2) + \lambda_2(-2\langle \triangle K^i - \triangle X^i, \triangle b_2^{i-1}\rangle$$
$$-\|\triangle K^i - \triangle X^i\|^2) + \lambda_3(-2\langle \triangle K^i - \triangle Z^i, \triangle b_3^{i-1}\rangle - \|\triangle K^i - \triangle Z^i\|^2)$$
$$\geq 2\lambda_1(\|\triangle K^i - \triangle M^i\|^2 + \langle \triangle M^i - \triangle M^{i-1}, \triangle K^i\rangle) + 2\lambda_2(\|\triangle K^i - \triangle X^i\|^2$$
$$+\langle \triangle X^i - \triangle X^{i-1}, \triangle K^i\rangle) + 2\lambda_3(\|\triangle K^i - \triangle Z^i\|^2 + \langle \triangle Z^i - \triangle Z^{i-1}, \triangle K^i\rangle)$$
$$-\lambda_1\|\triangle K^i - \triangle M^i\|^2 - \lambda_2\|\triangle K^i - \triangle X^i\|^2 - \lambda_3\|\triangle K^i - \triangle Z^i\|^2$$
$$= \lambda_1\|\triangle K^i - \triangle M^i\|^2 + \lambda_2\|\triangle K^i - \triangle X^i\|^2 + \lambda_3\|\triangle K^i - \triangle Z^i\|^2$$
$$+2\lambda_1\langle \triangle M^i - \triangle M^{i-1}, \triangle K^i\rangle + 2\lambda_2\langle \triangle X^i - \triangle X^{i-1}, \triangle K^i\rangle + 2\lambda_3\langle \triangle Z^i - \triangle Z^{i-1}, \triangle K^i\rangle.$$

Note that,

$$M^{i-1} = \arg\min_M h(M) + \frac{\lambda_1}{2}\|K^{i-1} - M + b_1^{i-2}\|_F^2,$$
$$X^{i-1} = \arg\min_X g(X) + \frac{\lambda_2}{2}\|K^{i-1} - X + b_2^{i-2}\|_F^2,$$
$$Z^{i-1} = \arg\min_Z \delta(Z) + \frac{\lambda_3}{2}\|K^{i-1} - Z + b_3^{i-2}\|_F^2,$$

so for any $M, X, Z \in \mathbb{R}^{m \times m}$,

$$h(M) - h(M^{i-1}) \geq (\nabla h(M^{i-1}))^T(M - M^{i-1}),$$

where $\nabla h(M^{i-1}) = -\lambda_1(-K^{i-1} + M^{i-1} - b_1^{i-2})$. Therefore,

$$h(M) - h(M^{i-1}) + \lambda_1\langle M^{i-1} - K^{i-1} - b_1^{i-2}, M - M^{i-1}\rangle \geq 0, \tag{37}$$

similar,

$$g(X) - g(X^{i-1}) + \lambda_2\langle X^{i-1} - K^{i-1} - b_2^{i-2}, X - X^{i-1}\rangle \quad \geq \quad 0, \tag{38}$$

$$\delta(Z) - \delta(Z^{i-1}) + \lambda_3\langle Z^{i-1} - K^{i-1} - b_3^{i-2}, Z - Z^{i-1}\rangle \quad \geq \quad 0. \tag{39}$$

Let $M = M^i, X = X^i, Z = Z^i$ in (37, 38, 39), then,

$$h(M^i) - h(M^{i-1}) + \lambda_1\langle M^{i-1} - K^{i-1} - b_1^{i-2}, M^i - M^{i-1}\rangle \quad \geq \quad 0,$$
$$g(X^i) - g(X^{i-1}) + \lambda_2\langle X^{i-1} - K^{i-1} - b_2^{i-2}, X^i - X^{i-1}\rangle \quad \geq \quad 0, \tag{40}$$
$$\delta(Z^i) - \delta(Z^{i-1}) + \lambda_3\langle Z^{i-1} - K^{i-1} - b_3^{i-2}, Z^i - Z^{i-1}\rangle \quad \geq \quad 0.$$

Let $M = M^{i-1}, X = X^{i-1}, Z = Z^{i-1}$ in (34), we have,

$$h(M^{i-1}) - h(M^i) + \lambda_1\langle M^i - K^i - b_1^{i-1}, M^{i-1} - M^i\rangle \quad \geq \quad 0,$$
$$g(X^{i-1}) - g(X^i) + \lambda_2\langle X^i - K^i - b_2^{i-1}, X^{i-1} - X^i\rangle \quad \geq \quad 0, \tag{41}$$
$$\delta(Z^{i-1}) - \delta(Z^i) + \lambda_3\langle Z^i - K^i - b_3^{i-1}, Z^{i-1} - Z^i\rangle \quad \geq \quad 0.$$

Summation of (40) and (41),

$$
\begin{aligned}
\langle b_1^{i-1} - b_1^{i-2}, M^i - M^{i-1}\rangle + \langle M^{i-1} - K^{i-1} - M^i + K^i, M^i - M^{i-1}\rangle &\geq 0, \\
\langle b_2^{i-1} - b_2^{i-2}, X^i - X^{i-1}\rangle + \langle X^{i-1} - K^{i-1} - X^i + K^i, X^i - X^{i-1}\rangle &\geq 0, \quad (42)\\
\langle b_3^{i-1} - b_3^{i-2}, Z^i - Z^{i-1}\rangle + \langle Z^{i-1} - K^{i-1} - Z^i + K^i, Z^i - Z^{i-1}\rangle &\geq 0.
\end{aligned}
$$

Also we know,

$$
\triangle K^i - \triangle K^{i-1} = K^i - K^{i-1}, \quad \triangle M^i - \triangle M^{i-1} = M^i - M^{i-1}, \quad \triangle X^i - \triangle X^{i-1} = X^i - X^{i-1},
$$
$$
\triangle Z^i - \triangle Z^{i-1} = Z^i - Z^{i-1}, \quad b_1^{i-1} - b_1^{i-2} = \triangle K^{i-1} - \triangle M^{i-1}, \quad b_2^{i-1} - b_2^{i-2} = \triangle K^{i-1} - \triangle X^{i-1},
$$
$$
b_3^{i-1} - b_3^{i-2} = \triangle K^{i-1} - \triangle Z^{i-1}.
$$

(42) can be written as

$$
\begin{aligned}
\langle \triangle K^{i-1} - \triangle M^{i-1}, \triangle M^i - \triangle M^{i-1}\rangle + \langle \triangle K^i - \triangle K^{i-1}, \triangle M^i - \triangle M^{i-1}\rangle &\geq \|\triangle M^i - \triangle M^{i-1}\|_F^2, \\
\langle \triangle K^{i-1} - \triangle X^{i-1}, \triangle X^i - \triangle X^{i-1}\rangle + \langle \triangle K^i - \triangle K^{i-1}, \triangle X^i - \triangle X^{i-1}\rangle &\geq \|\triangle X^i - \triangle X^{i-1}\|_F^2, \\
\langle \triangle K^{i-1} - \triangle Z^{i-1}, \triangle Z^i - \triangle Z^{i-1}\rangle + \langle \triangle K^i - \triangle K^{i-1}, \triangle Z^i - \triangle Z^{i-1}\rangle &\geq \|\triangle Z^i - \triangle Z^{i-1}\|_F^2.
\end{aligned}
$$

Because

$$
\begin{aligned}
\langle \triangle M^i - \triangle M^{i-1}, \triangle K^i\rangle &= \langle \triangle M^i - \triangle M^{i-1}, \triangle K^i - \triangle K^{i-1}\rangle \\
+ \langle \triangle M^i - \triangle M^{i-1}, \triangle K^{i-1} - \triangle M^{i-1}\rangle &+ \langle \triangle M^i - \triangle M^{i-1}, \triangle M^{i-1}\rangle, \\
\langle \triangle X^i - \triangle X^{i-1}, \triangle K^i\rangle &= \langle \triangle X^i - \triangle X^{i-1}, \triangle K^i - \triangle K^{i-1}\rangle \\
+ \langle \triangle X^i - \triangle X^{i-1}, \triangle K^{i-1} - \triangle X^{i-1}\rangle &+ \langle \triangle X^i - \triangle X^{i-1}, \triangle X^{i-1}\rangle, \\
\langle \triangle Z^i - \triangle Z^{i-1}, \triangle K^i\rangle &= \langle \triangle Z^i - \triangle Z^{i-1}, \triangle K^i - \triangle K^{i-1}\rangle \\
+ \langle \triangle Z^i - \triangle Z^{i-1}, \triangle K^{i-1} - \triangle Z^{i-1}\rangle &+ \langle \triangle Z^i - \triangle Z^{i-1}, \triangle Z^{i-1}\rangle,
\end{aligned}
$$

then

$$
\begin{aligned}
\langle \triangle M^i - \triangle M^{i-1}, \triangle K^i\rangle &\geq \|\triangle M^i - \triangle M^{i-1}\|_F^2 + \langle \triangle M^{i-1}, \triangle M^i - \triangle M^{i-1}\rangle, \\
\langle \triangle X^i - \triangle X^{i-1}, \triangle K^i\rangle &\geq \|\triangle X^i - \triangle X^{i-1}\|_F^2 + \langle \triangle X^{i-1}, \triangle X^i - \triangle X^{i-1}\rangle, \\
\langle \triangle Z^i - \triangle Z^{i-1}, \triangle K^i\rangle &\geq \|\triangle Z^i - \triangle Z^{i-1}\|_F^2 + \langle \triangle Z^{i-1}, \triangle Z^i - \triangle Z^{i-1}\rangle.
\end{aligned}
$$

So, (36) can be written as

$$
\begin{aligned}
\lambda_1(\|\triangle b_1^{i-1}\|_F^2 - \|\triangle b_1^i\|_F^2) &+ \lambda_2(\|\triangle b_2^{i-1}\|_F^2 - \|\triangle b_2^i\|_F^2) + \lambda_3(\|\triangle b_3^{i-1}\|_F^2 - \|\triangle b_3^i\|_F^2) \\
\geq \lambda_1\|\triangle K^i - \triangle M^i\|_F^2 &+ \lambda_2\|\triangle K^i - \triangle X^i\|_F^2 + \lambda_3\|\triangle K^i - \triangle Z^i\|_F^2 \\
&+ 2\lambda_1\|\triangle M^i - \triangle M^{i-1}\|_F^2 + 2\lambda_1\langle \triangle M^{i-1}, \triangle M^i - \triangle M^{i-1}\rangle \\
&+ 2\lambda_2\|\triangle X^i - \triangle X^{i-1}\|_F^2 + 2\lambda_2\langle \triangle X^{i-1}, \triangle X^i - \triangle X^{i-1}\rangle \\
&+ 2\lambda_3\|\triangle Z^i - \triangle Z^{i-1}\|_F^2 + 2\lambda_3\langle \triangle Z^{i-1}, \triangle Z^i - \triangle Z^{i-1}\rangle \\
= \lambda_1\|\triangle K^i - \triangle M^i\|_F^2 &+ \lambda_1\|\triangle M^i - \triangle M^{i-1}\|_F^2 + \lambda_1(\|\triangle M^i\|_F^2 - \|\triangle M^{i-1}\|_F^2) \\
+ \lambda_2\|\triangle K^i - \triangle X^i\|_F^2 &+ \lambda_2\|\triangle X^i - \triangle X^{i-1}\|_F^2 + \lambda_2(\|\triangle X^i\|_F^2 - \|\triangle X^{i-1}\|_F^2) \\
+ \lambda_3\|\triangle K^i - \triangle Z^i\|_F^2 &+ \lambda_3\|\triangle Z^i - \triangle Z^{i-1}\|_F^2 + \lambda_3(\|\triangle Z^i\|_F^2 - \|\triangle Z^{i-1}\|_F^2).
\end{aligned}
$$

Therefore,

$$(\lambda_1\|\triangle b_1^{i-1}\|_F^2 + \lambda_2\|\triangle b_2^{i-1}\|_F^2 + \|\triangle b_3^{i-1}\|_F^2 + \lambda_1\|\triangle M^{i-1}\|_F^2 + \lambda_2\|\triangle X^{i-1}\|_F^2 + \lambda_3\|\triangle Z^{i-1}\|_F^2)$$
$$-(\lambda_1\|\triangle b_1^i\|_F^2 + \lambda_2\|\triangle b_2^i\|_F^2 + \|\triangle b_3^i\|_F^2 + \lambda_1\|\triangle M^i\|_F^2 + \lambda_2\|\triangle X^i\|_F^2 + \lambda_3\|\triangle Z^i\|_F^2)$$
$$\geq \lambda_1(\|\triangle K^i - \triangle M^i\|_F^2 + \|\triangle M^i - \triangle M^{i-1}\|_F^2) + \lambda_2(\|\triangle K^i - \triangle X^i\|_F^2 + \|\triangle X^i - \triangle X^{i-1}\|_F^2)$$
$$+\lambda_3(\|\triangle K^i - \triangle Z^i\|_F^2 + \|\triangle Z^i - \triangle Z^{i-1}\|_F^2) \geq 0.$$

The sequence $\{\lambda_1\|\triangle b_1^i\|_F^2 + \lambda_2\|\triangle b_2^i\|_F^2 + \lambda_3\|\triangle b_3^i\|_F^2 + \lambda_1\|\triangle M^i\|_F^2 + \lambda_2\|\triangle X^i\|_F^2 + \lambda_3\|\triangle Z^i\|_F^2\}$ is non-increasing and convergent.

Also $\{K^i\},\{M^i\},\{X^i\},\{Z^i\},\{b_1^i\},\{b_2^i\},\{b_3^i\}$ are bounded and they have limit points. We get that $\lim_{i\to\infty}\|K^i - M^i\|_F^2 = 0$, $\lim_{i\to\infty}\|K^i - X^i\|_F^2$, $\lim_{i\to\infty}\|K^i - Z^i\|_F^2 = 0$.

Therefore, let $(\tilde{K},\tilde{M},\tilde{X},\tilde{Z},\tilde{b}_1,\tilde{b}_2,\tilde{b}_3)$ be limit point, i.e., there exists subsequence such that

$$\lim_{j\to\infty}(K^{i_j},M^{i_j},X^{i_j},Z^{i_j},b_1^{i_j},b_2^{i_j},b_3^{i_j}) = (\tilde{K},\tilde{M},\tilde{X},\tilde{Z},\tilde{b}_1,\tilde{b}_2,\tilde{b}_3).$$

In the following, we prove that $(\tilde{K},\tilde{M},\tilde{X},\tilde{Z})$ is minimum, i.e.,

$$f(\tilde{K}) + h(\tilde{M}) + g(\tilde{X}) + \delta(\tilde{Z}) = f(K^*) + h(M^*) + g(X^*) + \delta(Z^*).$$

Note that $(K^*,M^*,X^*,Z^*)$ is a saddle point, from lemma (3), $K^* = M^* = X^* = Z^*$,

$$f(K^*) + h(M^*) + g(X^*) + \delta(Z^*) \leq f(K^{i_j}) + h(M^{i_j}) + g(X^{i_j}) + \delta(Z^{i_j})$$
$$+\tfrac{\lambda_1}{2}\|K^{i_j} - M^{i_j}\|_F^2 + \tfrac{\lambda_2}{2}\|K^{i_j} - X^{i_j}\|_F^2 + \tfrac{\lambda_3}{2}\|K^{i_j} - Z^{i_j}\|_F^2$$
$$+\lambda_1\langle K^{i_j} - M^{i_j}, b_1^*\rangle + \lambda_2\langle K^{i_j} - X^{i_j}, b_2^*\rangle + \lambda_3\langle K^{i_j} - Z^{i_j}, b_3^*\rangle,$$

$j \to \infty$, we get that

$$f(K^*) + h(M^*) + g(X^*) + \delta(Z^*) \leq f(\tilde{K}) + h(\tilde{M}) + g(\tilde{X}) + \delta(\tilde{Z}). \tag{43}$$

On the other hand, let $K = K^*, M = M^*, X = X^*, Z = Z^*$ in equation (33) and (34), we get

$$f(K^*) + h(M^*) + g(X^*) + \delta(Z^*) \geq f(K^{i_j}) + h(M^{i_j}) + g(X^{i_j}) + \delta(Z^{i_j})$$
$$-\lambda_1\langle K^* - K^{i_j}, K^{i_j} - M^{i_j-1} + b_1^{i_j-1}\rangle - \lambda_2\langle K^* - K^{i_j}, K^{i_j} - X^{i_j-1} + b_2^{i_j-1}\rangle$$
$$-\lambda_3\langle K^* - K^{i_j}, K^{i_j} - Z^{i_j-1} + b_3^{i_j-1}\rangle - \lambda_1\langle M^* - M^{i_j}, M^{i_j} - K^{i_j} - b_1^{i_j-1}\rangle$$
$$-\lambda_2\langle X^* - X^{i_j}, X^{i_j} - K^{i_j} - b_2^{i_j-1}\rangle - \lambda_3\langle Z^* - Z^{i_j}, Z^{i_j} - K^{i_j} - b_3^{i_j-1}\rangle$$
$$= f(K^{i_j}) + h(M^{i_j}) + g(X^{i_j}) + \delta(Z^{i_j}) - \lambda_1\langle M^{i_j} - K^{i_j}, b_1^{i_j-1}\rangle$$
$$-\lambda_1\langle K^* - K^{i_j}, K^{i_j} - M^{i_j-1}\rangle - \lambda_1\langle M^* - M^{i_j}, M^{i_j} - K^{i_j}\rangle$$
$$-\lambda_2\langle X^{i_j} - K^{i_j}, b_2^{i_j-1}\rangle - \lambda_2\langle K^* - K^{i_j}, K^{i_j} - X^{i_j-1}\rangle$$
$$-\lambda_2\langle X^* - X^{i_j}, X^{i_j} - K^{i_j}\rangle - \lambda_3\langle Z^{i_j} - K^{i_j}, b_3^{i_j-1}\rangle$$
$$-\lambda_3\langle K^* - K^{i_j}, K^{i_j} - Z^{i_j-1}\rangle - \lambda_3\langle Z^* - Z^{i_j}, Z^{i_j} - K^{i_j}\rangle,$$

let $j \to \infty$, then

$$f(K^*) + h(M^*) + g(X^*) + \delta(Z^*) \geq f(\tilde{K}) + h(\tilde{M}) + g(\tilde{X}) + \delta(\tilde{Z}). \tag{44}$$

Combine with (43),

$$f(K^*) + h(M^*) + g(X^*) + \delta(Z^*) = f(\tilde{K}) + h(\tilde{M}) + g(\tilde{X}) + \delta(\tilde{Z}).$$

Hence, the limit point is minimum of (28). $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

# Acknowledgment

# References

[1] S. Choi, "Algorithms for orthogonal nonnegative matrix factorization," in *Neural Networks, 2008. IJCNN 2008.(IEEE World Congress on Computational Intelligence). IEEE International Joint Conference on*, pp. 1828–1832, IEEE, 2008.

[2] F. Pompili, N. Gillis, P.-A. Absil, and F. Glineur, "Two algorithms for orthogonal nonnegative matrix factorization with application to clustering," *Neurocomputing*, vol. 141, pp. 15–25, 2014.

[3] C. Ding, T. Li, W. Peng, and H. Park, "Orthogonal nonnegative matrix t-factorizations for clustering," in *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 126–135, ACM, 2006.

[4] J. Pan and M. K. Ng, "Orthogonal nonnegative matrix factorization by sparsity and nuclear norm optimization," *SIAM Journal on Matrix Analysis and Applications*, vol. 39, no. 2, pp. 856–875, 2018.

[5] L. De Lathauwer, B. De Moor, and J. Vandewalle, "On the best rank-1 and rank-(r 1, r 2,..., rn) approximation of higher-order tensors," *SIAM journal on Matrix Analysis and Applications*, vol. 21, no. 4, pp. 1324–1342, 2000.

[6] L. De Lathauwer, B. De Moor, and J. Vandewalle, "A multilinear singular value decomposition," *SIAM journal on Matrix Analysis and Applications*, vol. 21, no. 4, pp. 1253–1278, 2000.

[7] Y.-D. Kim and S. Choi, "Nonnegative tucker decomposition," in *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1–8, IEEE, 2007.

[8] R. Lai, J. Lu, and S. Osher, "Density matrix minimization with $\ell_1$ regularization," *Communications in Mathematical Sciences*, vol. 13, no. 8, pp. 2097–2117, 2015.

[9] H. Qi and D. Sun, "A quadratically convergent newton method for computing the nearest correlation matrix," *SIAM Journal on Matrix Analysis and Applications*, vol. 28, no. 2, pp. 360–385, 2006.

[10] P. Tseng, "Merit functions for semi-definite complementarity problems," *Mathematical Programming*, vol. 83, no. 2, pp. 159–186, 1998.

[11] D. Gabay and B. Mercier, "A dual algorithm for the solution of nonlinear variational problems via finite element approximation," *Computers & Mathematics with Applications*, vol. 2, no. 1, pp. 17–40, 1976.

[12] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, *et al.*, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends® in Machine learning*, vol. 3, no. 1, pp. 1–122, 2011.

[13] J. A. Hartigan and J. Hartigan, *Clustering algorithms*, vol. 209. Wiley New York, 1975.

[14] U. Von Luxburg, "A tutorial on spectral clustering," *Statistics and computing*, vol. 17, no. 4, pp. 395–416, 2007.

[15] J. MacQueen *et al.*, "Some methods for classification and analysis of multivariate observations," in *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, vol. 1, pp. 281–297, Oakland, CA, USA., 1967.

[16] H. Kim and H. Park, "Sparse non-negative matrix factorizations via alternating non-negativity-constrained least squares for microarray data analysis," *Bioinformatics*, vol. 23, no. 12, pp. 1495–1502, 2007.

[17] D. D. Lee and H. S. Seung, "Learning the parts of objects by non-negative matrix factorization," *Nature*, vol. 401, no. 6755, p. 788, 1999.

[18] M. W. Berry, M. Browne, A. N. Langville, V. P. Pauca, and R. J. Plemmons, "Algorithms and applications for approximate nonnegative matrix factorization," *Computational statistics & data analysis*, vol. 52, no. 1, pp. 155–173, 2007.

[19] F. S. Samaria and A. C. Harter, "Parameterisation of a stochastic model for human face identification," in *Applications of Computer Vision, 1994., Proceedings of the Second IEEE Workshop on*, pp. 138–142, IEEE, 1994.

[20] X. Li, M. K. Ng, G. Cong, Y. Ye, and Q. Wu, "Mr-ntd: manifold regularization nonnegative tucker decomposition for tensor data dimension reduction and representation," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 28, no. 8, pp. 1787–1800, 2017.

[21] A. H. Phan and A. Cichocki, "Tensor decompositions for feature extraction and classification of high dimensional datasets," *Nonlinear theory and its applications, IEICE*, vol. 1, no. 1, pp. 37–68, 2010.

[22] T.-K. Kim and R. Cipolla, "Canonical correlation analysis of video volume tensors for action categorization and detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 8, pp. 1415–1428, 2009.

[23] C. Schuldt, I. Laptev, and B. Caputo, "Recognizing human actions: a local svm approach," in *Pattern Recognition, 2004. ICPR 2004. Proceedings of the 17th International Conference on*, vol. 3, pp. 32–36, IEEE, 2004.

[24] F. Zhu, Y. Wang, B. Fan, S. Xiang, G. Meng, and C. Pan, "Spectral unmixing via data-guided sparsity," *IEEE Transactions on Image Processing*, vol. 23, no. 12, pp. 5412–5427, 2014.