

Joint Multiple Dictionary Learning for Tensor Sparse Coding

Yifan Fu, Junbin Gao

School of Computing and Mathematics
Charles Sturt University
Bathurst, NSW 2795, Australia
Email: {yfu, jbgao}@csu.edu.au

Yanfeng Sun

Beijing Municipal Key Lab of Multimedia
and Intelligent Software Technology
Beijing University of Technology
Beijing 100124, China
Email: yfsun@bjut.edu.cn

Xia Hong

School of Systems Engineering
University of Reading
Reading, RG6 6AY, UK
Email: x.hong@reading.ac.uk

Abstract—Traditional dictionary learning algorithms are used for finding a sparse representation on high dimensional data by transforming samples into a one-dimensional (1D) vector. This 1D model loses the inherent spatial structure property of data. An alternative solution is to employ Tensor Decomposition for dictionary learning on their original structural form—a tensor—by learning multiple dictionaries along each mode and the corresponding sparse representation in respect to the Kronecker product of these dictionaries. To learn tensor dictionaries along each mode, all the existing methods update each dictionary iteratively in an alternating manner. Because atoms from each mode dictionary jointly make contributions to the sparsity of tensor, existing works ignore atoms correlations between different mode dictionaries by treating each mode dictionary independently. In this paper, we propose a joint multiple dictionary learning method for tensor sparse coding, which explores atom correlations for sparse representation and updates multiple atoms from each mode dictionary simultaneously. In this algorithm, the Frequent-Pattern Tree (FP-tree) mining algorithm is employed to exploit frequent atom patterns in the sparse representation. Inspired by the idea of K-SVD, we develop a new dictionary update method that jointly updates elements in each pattern. Experimental results demonstrate our method outperforms other tensor based dictionary learning algorithms.

I. INTRODUCTION

Sparse Coding (SC) has been widely applied in numerous signal processing tasks, such as imaging denoising [1], texture synthesis [2] and image classification [3]. Sparse modeling aims at learning a dictionary so that each sample can be represented by a few atoms of the learned dictionary. Several algorithms have been developed for this task, e.g. the K-SVD [4] and the method of optimal directions (MOD) [5]. However, when dealing with multi-dimensional signals, all the previous sparse models reshape each input signal into a 1D vector. This kind of reshaping breaks the local correlation inherent inside the signal. Thus, it becomes highly desirable to develop algorithms which are able to fully make use of the local correlation inside high-dimensional data, such as tensor data.

Tensor decomposition [6] (TD) has attracted attention for processing multi-dimensional data. PARAFAC [6] and TUCKER [6] decompositions are two classical algorithms. PARAFAC decomposes a tensor as a sum of k rank-1 tensor while TUCKER factorizes a tensor into a set of matrices and one small core tensor. However, these methods do not

explicitly enforce sparsity constraint in the low dimensional representation.

Recently, researchers resort to combining SC with TD by introducing additional constraints to the models with the aim of learning sparse representations of tensors. Both non-negativity and sparsity have been used in two classical decompositions to accomplish the goal. Both non-negative versions of PARAFAC and TUCKER decompositions with multiplicative updates have been proposed in [7], [8] and [9], [10]. In TUCKER model, the sparsity over the core tensor is achieved by smoothing matrices along each mode [10], l_1 norm penalization [11] or a tensor dictionary learning algorithm [12].

In the case of dictionary learning models, Caiafa and Cichocki [13] open the discussion of sparse representation of tensor data using Kronecker bases in which the size of the core tensor is much higher than the input tensor. Two models are proposed in [13]: (1) Kronecker- Orthogonal Matching Pursuit (OMP) algorithm for multiway sparsity in which the sparse non-zero coefficients could be distributed randomly in the core tensor; (2) N-way Block OMP (N-BOMP) for multiway block sparsity in which the non-zero entries of the core tensor form blockwise structure. The similar ideas are used in 2 dimensional dictionaries for image processing [14]. In [12], a tensor dictionary model based on sparse TUCKER decomposition is proposed, in which sparse constraint over the core tensor is achieved by N-OMP, and the n -mode dictionary is learnt in an alternative minimization manner using gradient descent.

However, all existing dictionary learning algorithms update each mode dictionary by fixing all the other mode dictionaries iteratively in an alternating manner. The main issue is that the atoms from each mode dictionary jointly make contributions to the presentation of tensors, while current dictionary learning solutions ignore atoms correlations between different mode dictionaries. To this end, we propose a joint multiple dictionary learning method for tensor sparse coding (TSC-JMDL), which explores atom correlations for sparse representation and updates multiple atoms from each mode dictionary simultaneously. Our main contributions are as follows:

- Unlike previous tensor dictionary learning algorithms update all the non-zero atom entries for each mode, our model employs FP-tree to find the frequent atom patterns formed by a sequence of non-zero atom

entries in the sparse core tensor, and updates only frequent patterns by a tensor extended version of K-SVD algorithm.

- Existing tensor dictionary learning models update each mode dictionary independently, without considering atom correlations for sparse representation. In contrast, our model jointly updates all the elements in a frequent atom pattern simultaneously.

The remainder of the paper is organized as follows. Section II introduces the notations and the problem formulation. The algorithm details are introduced in Section III, followed by experimental results in Section IV. We conclude the paper in Section V.

II. NOTATIONS AND PROBLEM FORMULATION

A. Definition and Notations

A *tensor* is a multidimensional array. The *order* of a tensor is the number of dimensions, also known as ways or modes. For example, $\mathcal{Y} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ is an N -way tensor, where $I_n (1 \leq n \leq N)$ are the dimensions of each mode. The element indexed by (i_1, i_2, \dots, i_N) in an N -way tensor is denoted by y_{i_1, i_2, \dots, i_N} . In particular, a vector (1-way tensor) is denoted by a boldface lower-case letter, i.e. $\mathbf{y} \in \mathbb{R}^I$ and a matrix (2-way tensor) is denoted by a bold uppercase letter, i.e. $\mathbf{Y} \in \mathbb{R}^{I \times M}$. The i -th entry of a vector \mathbf{y} is denoted as y_i , and the element at (i, j) of a matrix \mathbf{Y} is denoted as y_{ij} . Thereafter, we will introduce some tensor fundamentals and definitions.

Definition 1 (Tensor Matricization): Matricization is the operation of rearranging the entries of a tensor so that it can be represented as a matrix. Let $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ be a tensor of order- N , the mode- n matricization of \mathcal{X} reorders the mode- n vectors to be columns of the resulting matrix, denoted by $\mathbf{X}_{(n)} \in \mathbb{R}^{I_n \times (I_{n+1} I_{n+2} \dots I_N I_1 I_2 \dots I_{n-1})}$.

Definition 2 (Rank-One Tensor): An N -way tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ is rank one if it can be equal to the outer product of N vectors:

$$\mathcal{X} = \mathbf{a}^{(1)} \circ \mathbf{a}^{(2)} \circ \dots \circ \mathbf{a}^{(N)}. \quad (1)$$

Elementwise, we have

$$x_{i_1 i_2 \dots i_N} = a_{i_1}^{(1)} a_{i_2}^{(2)} \dots a_{i_N}^{(N)} \text{ for all } 1 \leq i_n \leq I_n. \quad (2)$$

Definition 3 (Kronecker Product): The Kronecker product of matrices $\mathbf{A} \in \mathbb{R}^{I \times J}$ and $\mathbf{B} \in \mathbb{R}^{P \times L}$, denoted by $\mathbf{A} \otimes \mathbf{B}$, is a matrix of size $(IP) \times (JL)$ defined by

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \dots & a_{1J}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \dots & a_{2J}\mathbf{B} \\ \vdots & \vdots & \ddots & \vdots \\ a_{I1}\mathbf{B} & a_{I2}\mathbf{B} & \dots & a_{IJ}\mathbf{B} \end{bmatrix} \quad (3)$$

Definition 4 (The n -mode Product): The n -mode product of a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ by a matrix $\mathbf{U} \in \mathbb{R}^{J \times I_n}$, denoted as $\mathcal{X} \times_n \mathbf{U}$, is a tensor with entries:

$$(\mathcal{X} \times_n \mathbf{U})_{i_1, \dots, i_{n-1}, j, i_{n+1}, \dots, i_N} = \sum_{i_n=1}^{I_n} x_{i_1 i_2 \dots i_N} u_{j i_n} \quad (4)$$

The n -mode product is also denoted by each mode- n vector multiplied by the matrix \mathbf{U} . Thus, it can be expressed in terms of tensor matricization as well:

$$\mathcal{Y} = \mathcal{X} \times_n \mathbf{U} \quad \Leftrightarrow \quad \mathbf{Y}_{(n)} = \mathbf{U} \mathbf{X}_{(n)} \quad (5)$$

Definition 5 (Tucker Decomposition): Given an order- N tensor \mathcal{Y} , its Tucker decomposition is an approximated tensor defined by,

$$\hat{\mathcal{Y}} \equiv \llbracket \mathcal{X}; \mathbf{U}_1, \dots, \mathbf{U}_N \rrbracket = \mathcal{X} \times_1 \mathbf{U}_1 \times_2 \dots \times_N \mathbf{U}_N \\ = \sum_{i_1=1}^{M_1} \sum_{i_2=1}^{M_2} \dots \sum_{i_N=1}^{M_N} x_{i_1 i_2 \dots i_N} \mathbf{u}_{i_1} \circ \mathbf{u}_{i_2} \dots \circ \mathbf{u}_{i_N} \quad (6)$$

where $\mathcal{X} \in \mathbb{R}^{M_1 \times M_2 \times \dots \times M_N}$ is called a core tensor, $\mathbf{U}_i \in \mathbb{R}^{I_i \times M_i} (1 \leq i \leq N)$ are the factor matrices and the symbol \circ represents the vector outer product.

Definition 6 (CP Decomposition): Given an N -way tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$, the CP tensor decomposition factorizes a tensor into a sum of component rank-one tensors, which is defined as

$$\hat{\mathcal{X}} \equiv \sum_{i=1}^M \mathbf{u}_i^{(1)} \circ \mathbf{u}_i^{(2)} \dots \circ \mathbf{u}_i^{(N)} \quad (7)$$

where M is a positive integer and $\mathbf{u}_i^{(j)} \in \mathbb{R}^{I_j}$ for $i = 1, \dots, M$.

Definition 7 (Multiway block-sparsity): A tensor $\mathcal{Y} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ is (S_1, S_2, \dots, S_N) -block sparse with respect to the factors $\mathbf{U}_n \in \mathbb{R}^{I_n \times M_n} (n = 1, 2, \dots, N)$ if it admits a TUCKER representation based only on few S_n selected columns of each factor ($S_n \leq M_n$), i.e. if $\mathbf{i}_n = [i_n^1, i_n^2, \dots, i_n^{S_n}]$ denotes a subset of indices for mode $n (n = 1, 2, \dots, N)$, then

$$\mathcal{Y} = \mathcal{X} \times_1 \mathbf{U}_1 \times_2 \dots \times_N \mathbf{U}_N \quad (8)$$

with $x_{i_1 i_2 \dots i_N} = 0 \forall (i_1 i_2 \dots i_N) \notin \mathbf{i}_1 \times \mathbf{i}_2 \times \dots \times \mathbf{i}_N$

B. Problem Formulation

Given a set of M N -order tensors, denoted by $\mathcal{Y} = \{\mathcal{Y}_1, \mathcal{Y}_2, \dots, \mathcal{Y}_M\}$, we aim to find sparse representation for each N -order tensor $\mathcal{Y}_m (1 \leq m \leq M)$ with regards to the Kronecker product of multiple dictionaries \mathbf{U}_n for $1 \leq n \leq N$. Each \mathbf{U}_n is a dictionary along a particular direction of structure. To learn multiple dictionaries, we consider TUCKER decomposition on \mathcal{Y}_m . Thus the problem of multiple dictionary learning is formulated as

$$\min_{\mathcal{X}_m \text{ sparse}, \mathbf{U}_n} \sum_{m=1}^M \|\mathcal{Y}_m - \mathcal{X}_m \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \dots \times_N \mathbf{U}_N\| \quad (9)$$

where \mathcal{X}_m is the sparse representation tensor for \mathcal{Y}_m , and the factor matrix \mathbf{U}_n is the dictionary at mode- n . In this paper, we suppose \mathcal{X}_m is multiple block sparse as described in [13].

Naturally, we can regard \mathcal{Y} an $(N + 1)$ -order tensor by stacking all the given tensors along the $(N + 1)$ -mode. Similarly \mathcal{X} denotes the $(N + 1)$ -order tensor consisting of $\{\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_M\}$. Then model (9) can be equivalently written as

$$\min_{\mathcal{X}_m \text{ sparse}, \mathbf{U}_n} \|\mathcal{Y} - \mathcal{X} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \dots \times_N \mathbf{U}_N \times_{N+1} \mathbf{I}\| \quad (10)$$

where \mathbf{I} is the identity matrix of order M .

III. JOINT MULTIPLE DICTIONARY LEARNING MODEL FOR TENSOR SPARSE CODING

Similar to standard dictionary learning algorithms, the proposed tensor dictionary learning algorithm is a two-stage iterative process: sparse coding and dictionary update. In this paper, we employ an iterative algorithm called the Block Coordinate Descent (BCD) [15] to solve the optimization problem (10) by fixing all the other model variables to solve one variable at a time alternatively. Firstly, the dictionaries \mathbf{U}_n are fixed, the sparse coefficients \mathcal{X}_m can be obtained by solving M independent sparse representation subproblems. That is, for each $m = 1, \dots, M$, \mathcal{X}_m is obtained by minimizing

$$\min_{\mathcal{X}_m \text{ sparse}} \|\mathcal{Y}_m - \mathcal{X}_m \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \cdots \times_N \mathbf{U}_N\| \quad (11)$$

with the N-BOMP algorithm.

Then in the second stage, given the sparse representations \mathcal{X}_m , the dictionaries (factors) corresponding to each mode are updated jointly using an atom pattern based dictionary update strategy, by solving the following problem

$$\min_{\mathbf{U}_n} \|\mathcal{Y} - \mathcal{X} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \cdots \times_N \mathbf{U}_N \times_{N+1} \mathbf{I}\| \quad (12)$$

where $1 \leq n \leq N$.

A. N-BOMP Algorithm

OMP iteratively refines a sparse representation by successively identifying one component at a time that yields the greatest improvement in quality until a desired sparsity level is reached or the approximation error is below a pre-defined threshold. A tensor extended version of OMP for multiway block-sparsity is motivated by the fact that, in the real world, the non-zero coefficients are likely to be grouped in blocks rather than evenly distributed. Accordingly, N-BOMP is proposed to find a (S_1, S_2, \dots, S_N) -block sparse representation of an N -mode tensor with respect to the factors $\mathbf{U}_n \in \mathbb{R}^{I_n \times M_n}$ ($n = 1, 2, \dots, N$). If we denote by $\mathbf{B}_n \in \mathbb{R}^{I_n \times S_n}$ the submatrices formed by the columns indicated by indices \mathbf{i}_n of the mode- n dictionary, i.e. $\mathbf{B}_n = \mathbf{U}_n(:, \mathbf{i}_n)$, then the approximation of the tensor can be written in the equivalent vector version of Eq. (8) in terms of Kronecker products of dictionaries, that is

$$\hat{\mathbf{y}} = (\mathbf{B}_N \otimes \mathbf{B}_{N-1} \otimes \cdots \otimes \mathbf{B}_1) \mathbf{x} \quad (13)$$

where $\hat{\mathbf{y}} \in \mathbb{R}^{I_1 I_2 \cdots I_N}$ is the vectorized version of original tensor \mathcal{X} by stacking all the columns of mode-1 tensor $\mathcal{Y}_{(1)}$ in a single vector, and $\mathbf{x} \in \mathbb{R}^{S_1 S_2 \cdots S_N}$ is the vectorized version of the N -mode tensor consisting only of non-zero entries. Here we assume that the size of the core tensor \mathcal{X} is not less than the size of \mathcal{Y} ($M_n \geq I_n$), because sparse coding model is formulated with overcomplete dictionaries. The N-BOMP algorithm is given in Algorithm 1.

Remark 1: N-BOMP not only optimizes the memory usage but also requires much less iterations (K_{\max} against $S = S_1 S_2 \cdots S_N$ in the classical OMP/Kronecker-OMP algorithm), where S is the number of non-zero entries within the core tensor \mathcal{X} . Besides, the N-BOMP algorithm complexity in terms of the number of entries $I_1 I_2 \cdots I_N$, is sublinear compared to a linear dependence of the standard OMP and the Kronecker-OMP algorithms.

Algorithm 1 Solving Problem (11) by N-BOMP

Require: data tensor $\mathcal{Y}_m \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_N}$, mode- n dictionaries $\{\mathbf{U}_i\} (1 \leq i \leq N)$ with $\mathbf{U}_i \in \mathbb{R}^{I_i \times M_i}$, the maximum number of non-zero entries K_{\max} , error threshold ϵ

Ensure: Sparse representation $\mathcal{Y}_m = \mathcal{X}_m \times_1 \mathbf{U}_1 \times_2 \cdots \times_N \mathbf{U}_N$ with $x_{i_1 i_2 \cdots i_N} = 0 \forall (i_1, i_2, \dots, i_N) \notin \mathbf{i}_1 \times \mathbf{i}_2 \times \cdots \times \mathbf{i}_N$ (with non-zero entries given by $\mathcal{X}_m(\mathbf{i}_1, \mathbf{i}_2, \dots, \mathbf{i}_N) = \mathcal{A}$)

- 1: $\mathbf{i}_n = [\emptyset] (n = 1, 2, \dots, N)$, $\mathcal{R} = \mathcal{Y}_m$, $\mathcal{X}_m = 0$, $k = 1$;
- 2: **while** $|\mathbf{i}_1| |\mathbf{i}_2| \cdots |\mathbf{i}_N| \leq K_{\max}$ and $\|\mathcal{R}\|_F > \epsilon$ **do**
- 3: $[\mathbf{i}_1^k \mathbf{i}_2^k \cdots \mathbf{i}_N^k] = \arg \max_{[\mathbf{i}_1 \mathbf{i}_2 \cdots \mathbf{i}_N]} |\mathcal{R} \times_1 \mathbf{U}_1^T(:, \mathbf{i}_1) \times_2 \cdots \times_N \mathbf{U}_N^T(:, \mathbf{i}_N)|$;
- 4: $\mathbf{i}_n = \mathbf{i}_n \cup [\mathbf{i}_n^k] (n = 1, 2, \dots, N)$, $\mathbf{B}_n = \mathbf{U}_n(:, \mathbf{i}_n)$;
- 5: $\mathbf{a} = \arg \min_{\mathbf{w}} \|(\mathbf{B}_N \otimes \mathbf{B}_{N-1} \otimes \cdots \otimes \mathbf{B}_1) \mathbf{w} - \mathbf{y}\|_2^2$;
- 6: $\mathcal{R} = \mathcal{Y}_m - \mathcal{A} \times_1 \mathbf{B}_1 \times_2 \mathbf{B}_2 \cdots \times_N \mathbf{B}_N$;
- 7: $k = k + 1$;
- 8: **end while**
- 9: **return** $\{\mathbf{i}_1, \mathbf{i}_2, \dots, \mathbf{i}_N\}, \mathcal{A}$;

B. Joint multiple dictionary update algorithm

Once the sparse core tensors \mathcal{X}_m ($1 \leq m \leq M$) are obtained, the tensor dictionaries are computed by solving the problem (12). Because non-zero atoms entries from each mode dictionary jointly make contributions to the presentation of tensors, we propose a joint multiple dictionary update algorithm based on the frequent atom patterns for sparse tensor representation.

1) *Mining Frequent Atom Patterns in Sparse Representation:* Each non-zero entry of the sparse tensor \mathcal{X}_m ($1 \leq m \leq M$) is mapped into an integer using a function ϕ , which is formulated as

$$\begin{aligned} \phi(i_1 i_2 \cdots i_N) &= (i_1 - 1) I_2 I_3 \cdots I_N \\ &\quad + (i_2 - 1) I_3 I_4 \cdots I_N \\ &\quad + \cdots + i_N \end{aligned} \quad (14)$$

Then the sparse representation of the N -mode tensor \mathcal{X}_m ($1 \leq m \leq M$) is denoted by an integer set \mathbf{t}_i consisting of indices of the non-zero entries of \mathcal{X}_m .

In order to uncover the atom correlations for sparse representation with a low time cost, we adopt FP-tree [16] to find the frequent patterns from M integer sets. The time efficiency of FP-Tree is achieved with three techniques: (1) FP-Tree avoids costly repeated data base scans by compressing a large database into a condensed smaller data structure; (2) a pattern-fragment growth method is employed to avoid the costly generation of a large number of candidate sets; and (3) it dramatically reduces the search space by decomposing the mining task into a set of smaller tasks for mining confined patterns in conditional database. Given a database $\mathbf{T} = \{\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_M\}$ and a minimum support threshold ξ , FP-Tree mines frequent patterns (i.e. the support of patterns is not less than ξ) by creating conditional (sub)pattern-bases (i.e., the subpattern-base under the condition of a frequent item's existence). Taking Figure 1 as an example, let the sparse representation database \mathbf{T} be the first two columns of Fig. 1(a), and the minimum support threshold be 3 (i.e., $\xi = 3$). Firstly, a scan of \mathbf{T} derives a list of frequent items in Fig. 1(b), (the number after

“:” indicates the support), in which items are in frequency-descending order. Since each path of a tree will follow this order, the frequent atoms in each sparse representation are listed in the ordering in the rightmost column of Fig. 1(a). Secondly, the FP-Tree is constructed by scanning the database \mathbf{T} again, together with the associated node-links pointing to the nodes with the same atom in the tree, as shown in Fig. 1(b). Finally, the frequent patterns related to each frequent atom are mined by traversing the FP-Tree once following corresponding node-links. The conditional pattern-base and conditional FP-Trees are generated during this process, as shown in Fig. 1(c). Hence, the frequent patterns related to a specific frequent atom are summarized in the second column of Fig. 1(c).

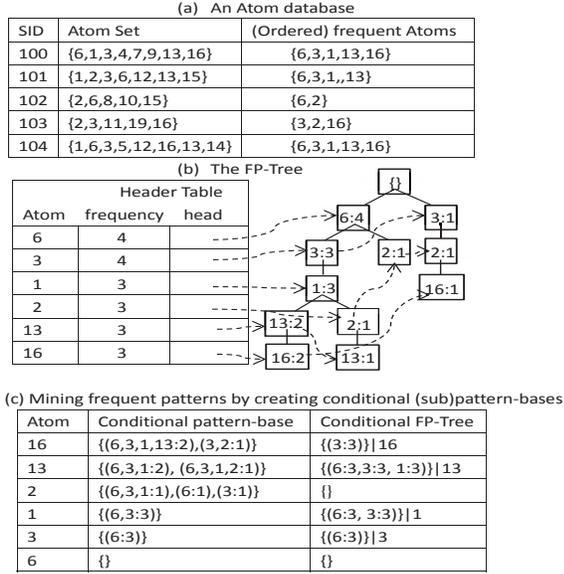


Fig. 1: A toy database illustrates the FP-Tree frequent pattern mining process

2) *Frequent Pattern based Dictionary Update*: After finding frequent patterns related to each frequent atom, i.e., the conditional FP-Tree regarding a specific atom, we successively update all frequent patterns by jointly updating all the elements in each frequent pattern at a time. First of all, a map function φ is used to convert an integer into an N -dimensional vector, which is denoted by

$$\varphi(b) = i_1 i_2 \dots i_N \quad (15)$$

where

$$m_1 = \frac{b}{I_2 I_3 \dots I_N} \quad (16)$$

$$i_1 = m_1 + 1$$

$$i_{n+1} = \frac{\text{mod}(m_n)}{I_{n+2} \dots I_N} + 1 \quad (1 \leq n \leq N-2) \quad (17)$$

$$i_N = \text{mod}\left(\frac{\text{mod}(m_{N-1})}{I_N}\right) \quad (18)$$

Using Eq. (16)-(18), each element in the frequent pattern is converted into a position (i_1, i_2, \dots, i_N) of a sparse tensor $\mathcal{X}_m (1 \leq m \leq M)$. Then we define ω_k as the group of indices

of sparse tensors \mathcal{X}_m that use all the elements in the frequent pattern p_k . Thus

$$\omega_k = \{m | 1 \leq m \leq M, \mathcal{X}_{p_k}^m \neq 0\} \quad (19)$$

where $p_k = \{u_1^{k_1} u_2^{k_1} \dots u_N^{k_1}, \dots, u_1^{k_n} u_2^{k_n} \dots u_N^{k_n}\}$, and $u_1^{k_j} u_2^{k_j} \dots u_N^{k_j} (1 \leq j \leq n)$ is the mapped position of sparse tensor in the frequent pattern p_k .

We now turn to the second stage of updating the dictionaries together with the non-zero coefficients. Assume both \mathcal{X} and $\mathbf{U}_i (1 \leq i \leq N)$, and we put in question only the dictionary atoms in each frequent pattern p_k , and the coefficients in \mathcal{X} that corresponding to it. Let \mathcal{X}_{p_k} and \mathcal{Y}_{p_k} be an $(N+1)$ -order tensor formed by stacking the subset ω_k of N -mode tensors \mathcal{Y}_m and \mathcal{X}_m along the $(N+1)$ mode, respectively; $\tilde{\mathcal{X}}_{p_k}$ and $\tilde{\mathcal{Y}}_{p_k}$ be the $(N+1)$ -order tensor formed by N -mode tensors \mathcal{Y}_m and \mathcal{X}_m not in the subset ω_k . Then the penalty term of the objective function (12) can be rewritten as

$$\begin{aligned} & \|\mathcal{Y} - \mathcal{X} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \dots \times_N \mathbf{U}_N \times_{N+1} \mathbf{I}\|_F^2 \\ &= \|\mathcal{Y} - \tilde{\mathcal{X}}_{p_k} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \dots \times_N \mathbf{U}_N \times_{N+1} \mathbf{I} \\ & \quad - \mathcal{X}_{p_k} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \dots \times_N \mathbf{U}_N \times_{N+1} \mathbf{I}\|_F^2 \\ &= \|\mathcal{E}_{p_k} - \mathcal{X}_{p_k} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \dots \times_N \mathbf{U}_N \times_{N+1} \mathbf{I}\|_F^2 \end{aligned} \quad (20)$$

where \mathcal{E}_{p_k} stands for the error for all the $|\omega_k|$ examples when the dictionary atoms in the frequent pattern p_k are removed. We obtain $\mathcal{E}_{p_k}^R$ by choosing only the $(N+1)$ -order indices in ω_k , as shown in Eq.(21).

$$\mathcal{E}_{p_k}^R = \mathcal{Y}(:, :, \dots, :, \omega_{p_k}) - \mathcal{X}_{(i_1, i_2, \dots, i_N, \omega_{p_k})} \times_1 \mathbf{U}_1 \dots \times_N \mathbf{U}_N \times_{N+1} \mathbf{I} \quad (21)$$

where $(i_1, i_2, \dots, i_N) \notin p_k$. We conduct a rank $(|p_k|, |p_k|, \dots, |p_k|)$ CP decomposition for $\mathcal{E}_{p_k}^R$, in which the CP factors for mode 1, 2, \dots , N of the j -th rank one tensor are regarded as the updated dictionary atoms of the j -th element in the frequent pattern p_k . The general process of joint multiple dictionary update algorithm is illustrated in Algorithm 2.

C. The complete Algorithm

After iteratively solving subproblem (11) and (12) until the maximum iterations are achieved or the iteration converges to stop, we finally obtain multiple overcomplete dictionaries along each mode of input tensors $\mathcal{Y}_m (1 \leq m \leq M)$ and corresponding sparse tensor representations $\mathcal{X}_m (1 \leq m \leq M)$. Algorithm 3 outlines the whole process of joint multiple dictionary learning algorithm for tensor sparse coding.

IV. EXPERIMENTAL RESULTS

In this section, we present a set of experimental results on both synthetic and real data sets with multi dimensional information. The intension of these experiments is to demonstrate our new method TSC-JMDL's superiority over the state-of-art dictionary learning methods in computation complexity, memory usage, and image denoising.

A. Baseline Methods

As our proposed method is closely related to tensor dictionary models based on sparse TUCKER decomposition, we implement two multiple based dictionary learning algorithms

Algorithm 2 Solving Problem (12) by Joint Multiple Dictionary Update Algorithm

Require: data tensor $\mathcal{Y} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N \times M}$, sparse representation $\mathcal{X} \in \mathbb{R}^{M_1 \times M_2 \times \dots \times M_N \times M}$ for $1 \leq i \leq M$, a minimum support threshold ξ

Ensure: Dictionaries from each mode $\mathcal{Y} = \mathcal{X} \times_1 \mathbf{U}_1 \times_2 \dots \times_N \mathbf{U}_N$

- 1: Initialize a non-zero entry database $\mathbf{T} = \{\}$
- 2: **for** $i = 1$ to M **do**
- 3: An integer set for the i -th slice of \mathcal{X} : $n_{\mathcal{X}_m} \leftarrow$ non-zero entries in \mathcal{X}_m ;
- 4: $\mathbf{t}_i \leftarrow \{\}$
- 5: **for** $j = 1$ to $|n_{\mathcal{X}_m}|$ **do**
- 6: $\mathbf{t}_i \leftarrow \mathbf{t}_i \cup \phi(n_{\mathcal{X}_m}(j))$;
- 7: **end for**
- 8: $\mathbf{T} \leftarrow \mathbf{T} \cup \mathbf{t}_i$;
- 9: **end for**
- 10: a frequent pattern list $\mathbf{P} = \{p_1, p_2, \dots, p_u\} \leftarrow$ Apply the FP-Tree algorithm to \mathbf{T} ;
- 11: **for** $k = 1$ to u **do**
- 12: **for** $j = 1$ to $|p_k|$ **do**
- 13: $u_1^{k_j} u_2^{k_j} \dots u_N^{k_j} \leftarrow \varphi(p_k(j))$;
- 14: $p_k(j) \leftarrow u_1^{k_j} u_2^{k_j} \dots u_N^{k_j}$
- 15: **end for**
- 16: $\omega_k \leftarrow \{i | 1 \leq i \leq M, \mathcal{X}_{p_k}^i \neq 0\}$;
- 17: Calculate the representation error tensor $\mathcal{E}_{p_k}^R$ with Eq.(21);
- 18: Do CP decomposition to the error tensor by a rank $(|p_k|, |p_k|, \dots, |p_k|)$ approximation:
 $\mathcal{E}_{p_k}^R \leftarrow \sum_{i=1}^{p_k} \mathbf{a}_{p_k}^{i(1)} \circ \dots \circ \mathbf{a}_{p_k}^{i(N)}$
- 19: **for** $j = 1$ to $|p_k|$ **do**
- 20: Update the dictionary atoms in the j -th element in p_k : $u_i^{k_j} = \mathbf{a}_{p_k}^{j(i)}$
- 21: **end for**
- 22: **end for**
- 23: **return** $\{\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_N\}$;

for tensor sparse representation as baseline methods. All use the N-BOMP algorithm for sparse representation learning. However, they make use of two different dictionary update methods. These two baselines are listed as follows:

- **Tensor MOD Approach (TMOD):** the original MOD algorithm is presented in [5] for dictionary learning on signal vectors. We apply this algorithm to tensor based multiple dictionary learning, which alternatively updates \mathbf{U}_n by fixing $\mathbf{U}_1, \dots, \mathbf{U}_{n-1}, \mathbf{U}_{n+1}, \dots, \mathbf{U}_N$ to minimizing (12)
- **Tensor K-SVD Approach (TKSVD):** Due to the atoms from each mode dictionary jointly make contributions to the tensor sparse representation, the K-SVD algorithm is unable to directly applied to tensor based dictionary learning by updating dictionary atom one by one. Instead, we check the coefficients $x_{i_1, i_2, \dots, i_N}^m$ at each position (i_1, i_2, \dots, i_N) of all the sparse tensors $\mathcal{X}_m (1 \leq m \leq M)$. For those $x_{i_1, i_2, \dots, i_N}^m \neq 0$, we gather the corresponding data tensor \mathcal{X}_m to generate an $(N + 1)$ -order tensor \mathcal{X}_s , and conduct a rank- $(1, 1, \dots, 1)$ CP decomposition. The CP factors for modes $1, 2, \dots, N$ are regarded as the updated dictio-

Algorithm 3 Joint multiple dictionary learning for tensor sparse coding

Require: data tensors $\mathcal{Y}_m \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N} (1 \leq m \leq M)$, the maximum number of non-zero entries K_{max} , error threshold ϵ , a minimum support threshold ξ

Ensure: Sparse representation corresponding the learnt multiple dictionaries along each structure mode $\mathcal{Y}_m = \mathcal{X}_m \times_1 \mathbf{U}_1 \times_2 \dots \times_N \mathbf{U}_N (1 \leq m \leq M)$

- 1: Initialize dictionaries along each mode $\mathbf{U}_j (1 \leq j \leq N)$;
- 2: **while** reach maximum iteration times or converge to stop **do**
- 3: Get the sparse representation tensors $\mathcal{X}_m (1 \leq m \leq M)$ by using N-BOMP;
- 4: Update dictionaries $\mathbf{U}_j (1 \leq j \leq N)$ by using TSC-JMDL;
- 5: **end while**
- 6: **return** Sparse representation $\mathcal{X}_m (1 \leq m \leq M)$ and Dictionaries for each mode $\mathbf{U}_j (1 \leq j \leq N)$

nary atoms.

We also consider some other sparse representation models proposed in [17], [18] and [14] to evaluate the performance in image denoising.

B. Synthetic Datasets

In this section, we evaluate TSC-JMDL against TMOD and TKSVD on the synthetic datasets. An $(N + 1)$ -order tensor of size $I_1 \times I_2 \dots \times I_N \times 100 (I_n = 10 \text{ for } n = 1, 2, \dots, N)$ is generated from the mode dictionaries of size $I_n \times M_n$ and the sparse core tensor of size $M_1 \times M_2 \dots \times M_N \times 100$ whose elements are obtained from Gaussian distributions, where $M_n = 2I_n$. The sparse core tensor has a fixed mode sparsity of $S_n = 5$.

1) *Time Cost with High Order Tensorial Data:* To show TSC-JMDL's time advantage of handling high order tensorial data over other baseline methods, we create 6 higher order tensors $\mathcal{X}_j (2 \leq j \leq 7)$ using above tensor construction method. Fig. 2 reports the results on TMOD, TKSVD and our new method TSC-JMDL. As the order of tensor increases, the time cost of TSC-JMDL and TMOD are significantly reduced compared with TKSVD. This is mainly because TKSVD jointly updates atoms, one for each mode at a time, the total number of these joint atoms is 20^{j-1} , which makes the dictionary update procedure particularly tedious. While TMOD simply updates each mode dictionary separately, the computation complexity of which is linear to the order of tensor. TSC-JMDL only jointly updates the frequent atom patterns in the sparse representation, which speeds up the dictionary update process. As TSC-JMDL needs extra time to find these frequent patterns, the time cost of TSC-JMDL is a little expensive than TMOD. To sum up, our method TSC-JMDL has a comparable time cost in the higher mode of tensor.

2) *Data Denoising with Different Noise Parameters:* We investigate the performance of TSC-JMDL for data denoising by comparing with TMOD and TKSVD. The test data 3-order tensor is generated using above the same tensor construction method. The Gaussian noise is artificially added with zero

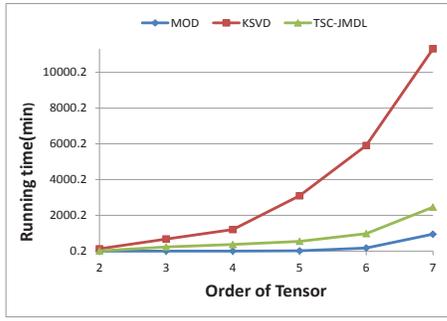


Fig. 2: Time comparison *w.r.t.* different orders of a tensor.

mean. Fig. 3 presents the comparison of peak signal-to-noise ratios (PSNR) of the denoising results with increasing the noise parameter σ from 5 to 100. It can be found that joint multiple dictionary update algorithms TKSVD and TSC-JMDL outperform independent dictionary update algorithm TMOD, which suggests that considering atoms correlations for sparse representation can boost performance in denoising. Moreover, the performance of TKSVD is marginally better than the proposed method TSC-JMDL. This is mainly because TSC-JMDL updates merely frequent atoms patterns rather than all the non-zero coefficients in the sparse tensors. However, our new method TSC-JMDL can achieve a comparable performance with much lower time cost than TKSVD.

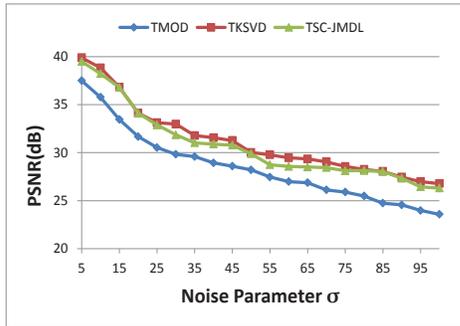


Fig. 3: The PSNR of denoising results with the variation of σ .

3) *Memory Usage with respect to Dictionaries:* Exemplified dictionaries used in a 3-order tensor (created using above tensor generation method) is illustrated in Fig. 4. The left two images presents the dictionaries U_1 and U_2 learnt by our method TSC-JMDL, in which each column is the atom along one direction of the data. The right image is the Kronecker product U of above two dictionaries, which fully represents the data spatial correlation. Note that the size of dictionaries are quite different. The classical OMP requires a large dictionary of size 100×400 , whereas N-BOMP employed in our scheme needs only two dictionaries of size 10×20 . Though only $1/100$ size of the Kronecker dictionary is used, our model can achieve better performance than the classical OMP algorithm. This is mainly because our model simultaneously selects factors along each mode on the original structured data, while 1D sparse model OMP chooses the factors on the reshaped 1D vector.

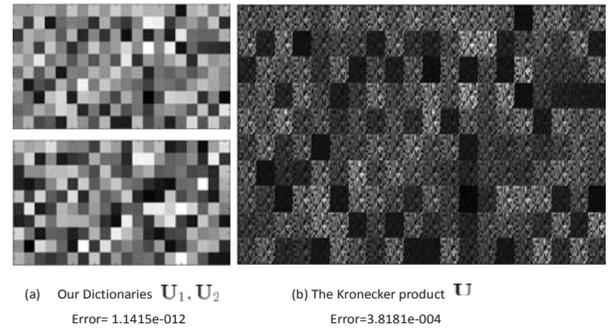


Fig. 4: Exemplified dictionaries in our 3-order tensor model.

C. Real Datasets

In this section, our proposed TSC-JMDL is used to solve the image denoising problem. The test images include 5 widely used images [18], commonly known as “Lena”, “Peppers”, “House”, “Barbara”, and “Boats”, which are resized uniformly into 512×512 . The white Gaussian noises are added at different standard deviations σ . In all the tests, image patch is of size 8×8 pixels, and the two dictionaries U_1 and U_2 are of size 8×16 . Then U is generated by using Kronecker product is of size 64×256 .

1) *Image Denoising with Different Noise Variances:* We study the performance of our TSC-JMDL in image denoising with those of 1D sparse coding models proposed in [17] and [18] and one 2D synthesis sparse model proposed in [14]. Table I shows the denoising results in terms of PSNR. We observe that our method outperforms the 2D synthesis sparse model [14] among all the data sets with different variations. This fact demonstrates that considering atom correlations for sparse representation can boost image denoising performance. Note that the size of dictionaries are different in generating the results. The methods in [17] and [18] require the dictionaries of size 64×256 , whereas our method only needs two dictionaries of size 8×16 , which is $1/64$ size of the dictionary used in the above two 1D sparse coding models. However, our model TSC-JMDL can achieve the best performance among all the four evaluated methods as denoted by the bolded numbers.

2) *Performance with the Same Size of Dictionaries:* We further study the performance of our method with the same size of dictionaries used in other baseline methods. Table II shows the denoising results of [14] and [18] using different sizes of dictionaries which are equal to the size used in our method. The noise parameter $\sigma = 5$, and the noise image of PSNR = $34.15dB$. Clearly, the larger the size of dictionaries, the higher PSNR results all the methods achieve. Besides, two dictionary sparse models (TSC-JMDL and 2D synthesis sparse model [14]) significantly outperforms the 1D sparse model [18] when same size of dictionary is used. Another interesting fact is that the denoising result of our proposed method TSC-JMDL always performs better than the 2D synthesis sparse model, which again demonstrates that taking atom correlations into consideration can improve denoising results.

A visual comparison is given in Fig. 5. It presents the denoising results of Lena, Peppers and House generated by 1D sparse model [18], 2D synthesis sparse model [14] and

TABLE I: Summary of the denoising PSNR results in [dB]. In each cell, the top row is the result of Portilla et al. [17], the second row is the result of Elad et al. [18]. They all use the dictionary of size 64×256 . The third row is the result of Qi et al. [14], and the bottom row is the result of the proposed method. They all use two dictionaries of size 8×16 .

| $\sigma \setminus PSNR$ | 2 \ 42.11 | 5 \ 34.15 | 10 \ 28.13 |
|-------------------------|--------------|--------------|--------------|
| Lena | 42.23 | 38.49 | 35.61 |
| | 43.58 | 38.6 | 35.47 |
| | 43.58 | 38.55 | 35.37 |
| | 47.01 | 42.12 | 39.96 |
| Barbara | 43.29 | 37.79 | 34.03 |
| | 43.67 | 38.08 | 34.42 |
| | 43.64 | 38.05 | 34.02 |
| | 44.65 | 39.77 | 36.37 |
| Peppers | 43 | 37.31 | 33.77 |
| | 43.33 | 37.78 | 34.28 |
| | 43.37 | 37.93 | 34.26 |
| | 46.57 | 41.98 | 37.02 |
| House | 44.07 | 38.65 | 35.35 |
| | 44.47 | 39.37 | 35.98 |
| | 44.38 | 39.14 | 35.59 |
| | 48.98 | 45.31 | 41.07 |
| Boats | 42.09 | 36.97 | 33.58 |
| | 43.14 | 38.08 | 33.64 |
| | 43.11 | 37.16 | 33.56 |
| | 43.42 | 38.01 | 34.39 |

TABLE II: Summary of the denoising PSNR results in [dB]. In each cell, the top row is the result of Elad et al. [18]. It uses the dictionaries of size 64×4 , 64×16 and 64×64 respectively. The middle row is the result of Qi et al. [14], and the bottom row is the result of the proposed method. They all use two dictionaries of size 8×16 , 8×64 and 8×256 respectively.

| <i>image</i> \ <i>DictSize(pixels)</i> | 256 | 1024 | 4096 |
|--|--------------|--------------|--------------|
| Lena | 30.09 | 35.81 | 38.20 |
| | 38.55 | 44.86 | 47.12 |
| | 42.12 | 47.25 | 49.94 |
| Barbara | 24.40 | 30.14 | 37.74 |
| | 38.05 | 41.36 | 46.76 |
| | 39.77 | 43.02 | 49.96 |
| Peppers | 25.73 | 32.23 | 37.29 |
| | 37.93 | 42.28 | 45.16 |
| | 41.98 | 45.98 | 48.79 |
| House | 29.35 | 36.15 | 39.34 |
| | 39.14 | 44.68 | 46.59 |
| | 45.31 | 47.99 | 49.27 |
| Boats | 27.08 | 33.04 | 37.14 |
| | 37.16 | 41.78 | 48.09 |
| | 38.01 | 43.03 | 49.78 |

our proposal model TSC-JMDL with the dictionary of the same sizes 64×4 and $2 \times 8 \times 16 = 256$, respectively. Obviously, our method TSC-JMDL provides much clearer reconstructed image than our two baseline methods.

Exemplified dictionaries of Lena used in our method is illustrated in Fig. 6. The 2D dictionaries U_1 , U_2 denote different dimensional features, and the Kronecker product U fully represents the image spatial correlations along each direction. Two dictionaries in our model are learnt by exploring the structure information along each mode and their correlations among all the possible modes, whereas the Kronecker product U is generated depending on one direction information on the reshaped image vector. Thus, our model always has better image denoising results than 1D sparse model.

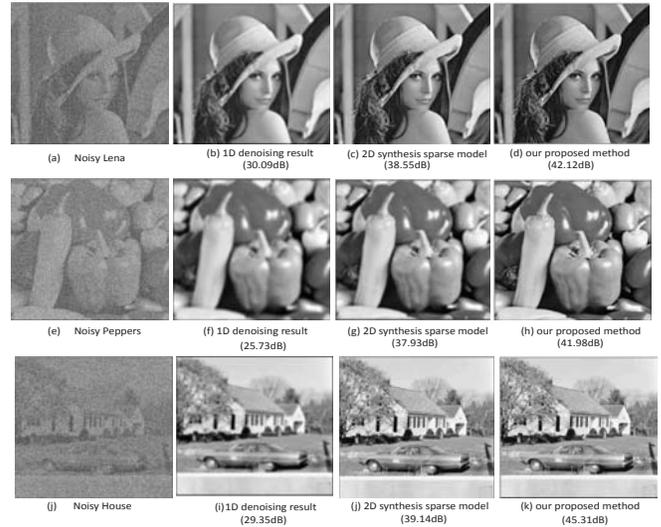


Fig. 5: The denoising results of Lena, Peppers and House by 1D sparse model [18], 2D synthesis sparse model [14] and our proposal model TSC-JMDL using the dictionaries of size 64×4 and 8×16 , respectively.

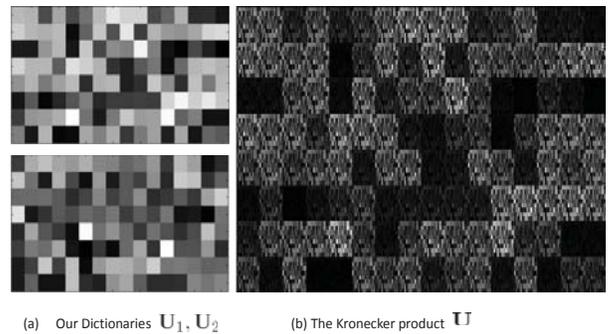


Fig. 6: Exemplified dictionaries of Lena in our model TSC-JMDL. Left two images show U_1 , U_2 in which each column is the atom of one directional signal of the image patch. The right image is the Kronecker dictionary U in which each square is an atom of size 8×8 .

V. CONCLUSION

In this paper, we propose a novel multiple dictionary learning algorithm for tensor sparse coding. While other existing tensor dictionary learning algorithms update each mode dictionary by fixing all the other mode dictionaries iteratively in an alternating manner, the proposed method fully makes use of the atom correlations among all the spatial modes inside a higher order tensorial data. Our model TSC-JMDL employs the FP-tree mining algorithm to exploit frequent atom patterns in the sparse representation, and then simultaneously updates multiple atoms in the frequent pattern using our joint dictionary update method. On the synthetic datasets, we show that our model can achieve a comparable performance with a lower computation gain and memory usage. Moreover, we also demonstrate its effectiveness in data denoising with different noise variances. On the real-world datasets, our method shows promising results in image denoising. Our model outperforms

both tensor dictionary learning methods and traditional 1D models, with different noise parameters and similar memory cost. The image reconstruction results clearly show the ability of our algorithm for maintaining the discerning features while retaining the image reconstruction.

ACKNOWLEDGMENT

This work is supported by the Australian Research Council (ARC) through Discovery Project Grant DP130100364.

REFERENCES

- [1] S. Li, "Non-negative sparse coding shrinkage for image denoising using normal inverse gaussian density model," *Image Vision Comput.*, vol. 26, no. 8, pp. 1137–1147, Aug. 2008. [Online]. Available: <http://dx.doi.org/10.1016/j.imavis.2007.12.006>
- [2] G. Peyré, "Sparse Modeling of Textures," *Journal of Mathematical Imaging and Vision*, vol. 34, no. 1, pp. 17–31, May 2009. [Online]. Available: <http://dx.doi.org/10.1007/s10851-008-0120-3>
- [3] J. Yang, K. Yu, Y. Gong, and T. Huang, "Linear spatial pyramid matching using sparse coding for image classification," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [4] M. Aharon, M. Elad, and A. Brucktein, "K-SVD: An algorithm for desdesign overcomplete dictionaries for sparse representation," *IEEE Trans. on Signal Processing*, vol. 54, no. 2, pp. 4311–4322, 2006.
- [5] K. Engan, S. Aase, and J. Husoy, "Method of optimal directions for frame design," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, vol. 1, 1999, pp. 2443–2446.
- [6] G. Kolda and B. Bader, "Tensor decompositions and applications," *SIAM Review*, vol. 51(3), pp. 455–500, 2009.
- [7] T. Hazan, S. Polak, and A. Shashua, "Sparse image coding using a 3d non-negative tensor factorization," in *Computer Vision, 2005. ICCV 2005. Tenth IEEE International Conference on*, vol. 1, 2005, pp. 50–57 Vol. 1.
- [8] E. Benetos and C. Kotropoulos, "Non-negative tensor factorization applied to music genre classification," *Audio, Speech, and Language Processing, IEEE Transactions on*, vol. 18, no. 8, pp. 1955–1967, 2010.
- [9] Y. Kim, A. Cichocki, and S. Choi, "Nonnegative tucker decomposition with alpha-divergence," in *ICASSP. IEEE*, 2008, pp. 1829–1832. [Online]. Available: <http://dblp.uni-trier.de/db/conf/icassp/icassp2008.html#KimCC08>
- [10] Y. Kim and S. Choi, "Nonnegative tucker decomposition," in *Computer Vision and Pattern Recognition, 2007. CVPR '07. IEEE Conference on*, 2007, pp. 1–8.
- [11] M. Mørup, L. Hansen, and S. M. Arnfred, "Algorithms for sparse nonnegative tucker decompositions," *Neural Comput.*, vol. 20, no. 8, pp. 2112–2131, Aug. 2008. [Online]. Available: <http://dx.doi.org/10.1162/neco.2008.11-06-407>
- [12] S. Zubair and W. Wang, "Tensor dictionary learning with sparse tucker decomposition," in *Digital Signal Processing (DSP), 2013 18th International Conference on*, 2013, pp. 1–6.
- [13] C. F. Caiafa and A. Cichocki, "Computing sparse representations of multidimensional signals using kronecker bases," *Neural Comput.*, vol. 25, no. 1, pp. 186–220, Jan. 2013. [Online]. Available: http://dx.doi.org/10.1162/NECO_a_00385
- [14] N. Qi, Y. Shi, X. Sun, J. Wang, and B. Yin, "Two dimensional synthesis sparse model," in *Multimedia and Expo (ICME), 2013 IEEE International Conference on*, 2013, pp. 1–6.
- [15] M. Blondel, K. Seki, and K. Uehara, "Block coordinate descent algorithms for large-scale sparse multiclass classification," *Machine Learning*, vol. 93, no. 1, pp. 31–52, 2013. [Online]. Available: <http://dx.doi.org/10.1007/s10994-013-5367-2>
- [16] J. Han, J. Pei, Y. Yin, and R. Mao, "Mining frequent patterns with candidate generation: A frequent-pattern tree approach," *Data Mining and Knowledge Discovery*, vol. 8, pp. 53–87, 2004.
- [17] J. Portilla, V. Strela, M. Wainwright, and E. Simoncelli, "Image denoising using scale mixtures of gaussian in the wavelet domain," *Image Processing, IEEE Transactions on*, vol. 12, pp. 1338–1351, 2003.
- [18] M. Elad, M. and AElad, "Image denoising via sparse and redundant representation over learned dictionaries," *Image Processing, IEEE Transactions on*, vol. 15, pp. 3736–3745, 2006.