

FacetCube: a general framework for non-negative tensor factorization

Yun Chi · Shenghuo Zhu

Received: 5 March 2010 / Revised: 23 July 2012 / Accepted: 29 July 2012
© Springer-Verlag London 2012

Abstract Non-negative tensor factorization (NTF) has been successfully used to extract significant characteristics from polyadic data, such as data in social networks. Because these polyadic data have multiple dimensions (e.g., the author, content, and timestamp of a blog post), NTF fits in naturally and extracts data characteristics *jointly* from different data dimensions. In the traditional NTF, all information comes from the observed data, and therefore, the end users have no control over the outcomes. However, in many applications very often, the end users have certain prior knowledge, such as the demographic information about individuals in a social network or a pre-constructed ontology on the contents and therefore prefer the data characteristics extracting by NTF being consistent with such prior knowledge. To allow users' prior knowledge to be naturally incorporated into NTF, in this paper, we present a general framework—FacetCube—that extends the standard NTF. The new framework allows the end users to control the factorization outputs at three different levels for each of the data dimensions. The proposed framework is intuitively appealing in that it has a close connection to the probabilistic generative models. In addition to introducing the framework, we provide an iterative algorithm for computing the optimal solution to the framework. We also develop an efficient implementation of the algorithm that consists of several techniques to make our framework scalable to large data sets. Extensive experimental studies on a paper citation data set and a blog data set demonstrate that our new framework is able to effectively incorporate users' prior knowledge, improves performance over the traditional NTF on the task of personalized recommendation, and is scalable to large data sets from real-life applications.

Keywords Non-negative tensor factorization · Polyadic data · Prior knowledge · Iterative algorithm · Sparse algorithm

Y. Chi (✉) · S. Zhu
NEC Laboratories America, Inc., 10080 North Wolfe Road,
SW3-350, Cupertino, CA 95014, USA
e-mail: ychi@nec-labs.com

1 Introduction

Data in many applications are *polyadic*, that is, they have multiple dimensions. Data in social networks are such an example—for example, data in the blogosphere may have people dimension (the author of a blog post), content dimension (the body of the post), and time dimension (the timestamp of the post). Documents in a digital library are another example—a scientific paper can be described by its authors, keywords, references, publication date, publication venue, etc. To analyze such polyadic data, a very important task is to extract significant characteristics from different data dimensions, where the extracted characteristics can be used either directly for data summarization and visualization or as features for further data analysis. The extracted characteristics can be in the form of, using the blog example, salient communities among bloggers, coherent topics in blog posts, and noteworthy temporal trends of these topics and communities. Because these data dimensions affect each other in a *joint* way, approaches that either analyze each data dimension independently or only consider pairwise relations between data dimensions are not able to accurately capture the data characteristics.

Recently, several multiple-dimensional tensor models have been proposed to capture the higher-order correlation (other than the second-order correlation) among various data dimensions. These tensor-based approaches can be categorized into two groups. Approaches in the first group [9, 16] decompose polyadic data by using higher-order *linear* decompositions such as higher-order singular value decomposition (HOSVD), which are extensions of the matrix singular value decomposition (SVD). On the other hand, approaches in the second group [7, 8, 27, 37] decompose polyadic data by using non-negative tensor factorizations (NTFs), which are extensions of the non-negative matrix factorization (NMF) [22]. Approaches based on NTFs decompose data into *additions* of non-negative components and therefore have many advantages over those based on linear decompositions. Such advantages include ease of interpretation of the extracted characteristics, close connection to the probabilistic models, and no enforcement on the orthogonality among different data characteristics. Because of these advantages, in this paper, we focus on the approaches that extract data characteristics by using NTF.

In an approach based on the standard NTF for extracting data characteristics, the extracted characteristics can be of arbitrary forms and end users do not have any control over them. Such an approach has some benefits—it is simple because it does not require any input other than the observed data. However, such a simple approach also has its weakness: end users have no channel to incorporate their prior knowledge into the process of characteristic extraction. We argue that users' prior knowledge can greatly benefit the extraction of meaningful data characteristics from different perspectives. Cognitively, human concepts such as the content topics usually only occupy a very low-dimensional subspace or manifold of the whole space (which usually is of much higher dimensions), and users' prior knowledge can provide guidance toward the appropriate subspace. Statistically, a troublesome issue in the standard NTF-based approach is the problem of overfitting, whereas users' input can alleviate this problem. Application-wise, there are many applications in which the end users already have certain domain knowledge, for example, a pre-constructed ontology of contents, and want to view the data through the lens of such domain knowledge. Therefore, it is beneficial to provide means to allow the end users to incorporate prior knowledge into the process of characteristic extraction.

In this paper, we propose a new framework that extends the standard NTF and allows end users to incorporate prior knowledge in the process of extracting characteristics from

polyadic data. We name our framework FacetCube.¹ Our main contributions are summarized as follows.

- We propose a novel and more general non-negative tensor factorization model for extracting data characteristics from polyadic data. In the new model, prior knowledge can be incorporated at three different levels into the factorization, and as a result, end users can control the process of characteristic extraction in different data dimensions at different levels. The new model takes the standard NTF as a special case.
- We show that the proposed model has a natural interpretation in a form of the probabilistic generative procedure. From this interpretation, we propose and justify the usage of a Dirichlet prior in the model parameter inference. We further show a technique of controlling the sparseness of the results by using a special Dirichlet prior.
- We develop a very efficient implementation of the algorithm that takes advantage of the sparseness of data, where the implementation has linear (per iteration) time complexity. The implementation contains a series of techniques that are general enough to be equally applicable to any algorithms with similar iterative computations. We further develop an efficient technique of fast computation when the application is to query about the top- K answers.

Our FacetCube framework is general in that it can handle polyadic data of arbitrary order. In this paper, we use two applications, personalized recommendation in a paper citation data set and data exploration in a blog data set, to demonstrate the effectiveness of our framework. Extensive experimental studies on these two data sets demonstrate that our new framework is able to effectively incorporate users' prior knowledge, improves performance over the standard NTF on the task of personalized recommendation, and is scalable to large data sets from real-life applications.

The rest of this paper is organized as the following. In the rest of this section, we survey related work. In Sect. 2, we provide background information. In Sect. 3 we introduce our FacetCube framework. In Sect. 4, we offer a probabilistic interpretation of our framework. In Sect. 5, we describe the techniques for our efficient implementation. We present experimental results in Sect. 6 and finally give conclusions in Sect. 7.

Related work

The majority of existing tensor factorization-based studies focus on *linear* tensor factorizations, where the non-negative constraint is not enforced. Two mostly used linear tensor factorizations are the PARAFAC model [16] and the Tucker model [33], where the latter has recently been generalized in [9] to the Higher Order Singular Value Decomposition. These models have been used for applications such as Web page ranking [32], cross-lingual information retrieval [4], sensor network analysis [31], etc. We refer interested readers to the comprehensive survey [20]. Although these linear tensor factorizations show some good properties in terms of minimizing certain losses in the Frobenius norm, they also have some weak points. For example, the negative values in the outcome of the factorization make it difficult to interpret the outcome in terms of probabilistic distributions. In addition, in contrast to NTF where the data are factorized into *additions* of components, linear tensor factorizations require *subtractions* among components which make them not intuitively appealing. Furthermore,

¹ *FacetCube* stands for “factorize data using NTF with co-ordinates being unconstrained, basis-constrained, or fixed”.

linear tensor factorizations usually require concepts in the outcome to be orthogonal, and this limits the scope of their applications (e.g., we seldom enforce different temporal trends to be orthogonal to each other). A special linear tensor factorization, the CANDELINC model [3], has certain similarity to our framework in that it allows the data characteristics to be inside a pre-defined linear subspace. However, it can be shown that for linear factorizations, such a goal can be trivially achieved by projecting the data into the subspace before applying the tensor factorization.

Our work is an extension to the existing NTFs, and the following are some related studies in this area. Hofmann proposed the probabilistic semantic analysis (PLSA) for extracting latent topics from documents [18]. Later, Gaussier and Goutte [15] and Ding et al. [11] showed the equivalence between the PLSA and the NMF algorithms [22]. Several studies have extended NMFs to NTFs and applied them to applications such as computer vision [30, 37], digital signal processing [14], information retrieval [7], etc. Both Tucker-typed NTF [7, 8, 27] and the PARAFAC-typed NTF [17, 30] have been explored in some of these studies. Our work extends existing NTF approaches by allowing users to incorporate prior knowledge in the factorization process and takes the above standard NTF approaches as special cases. In addition, in this work, we focus on the Tucker-typed NTF and the extension to the PARAFAC-typed NTF is straightforward. In addition, this paper is an extension of our previously published paper [5].

Some other related works include the following. Dhillon et al. [10] proposed a co-clustering framework using an information theoretic framework. Long et al. [26] proposed several optimization-based models for multi-type relational data. Wang et al. [34] studied the clustering problem from a viewpoint of learning a bi-stochastic data similarity matrix. These approaches, however, focus more on pairwise relations and do not capture higher-order correlations. Banerjee et al. [2] proposed a multi-way clustering framework that can handle multi-dimensional relations. However, the approach by Banerjee et al. is derived from the Bregman divergence, and it uses an optimization framework. In comparison, our new model has an underlying probabilistic interpretation. In addition, Chi et al. [6] proposed to incorporate certain regularization, such as the connectivity of subgraphs and the smoothness of temporal trends, in the factorization process. Lin et al. [23, 24] proposed a *FacetNet* framework for dynamic community detections. These two studies turned out to solve non-negative *matrix* factorization problems instead of using non-negative *tensor* factorizations. Peng [28] and Chi et al. [7] both pointed out the connection between NTF and tensorial PLSA. However, the model proposed in this paper takes prior knowledge into consideration and as a result, the probabilistic model is slightly different.

Recently, there have been some new models, such as the BPTF model in [35] and the Multi-HDP model in [29] that chose to use Bayesian treatments. Usually in a Bayesian model, the inference could be very expensive. However, it is interesting to compare the performance of our framework with these Bayesian frameworks, and this task is one of our future directions. Lin et al. [25] recently proposed the MetaFac model based on NTF, where the temporal correlation is considered in the factorization. However, in [25], the prior knowledge is used as a regularization term, and so is not enforced. In the field of information retrieval, there are many other studies on how to incorporate prior knowledge into an existing model. For example, Aussenac-Gilles and Mothe [1] proposed to use domain-specific ontology to help explore documents. For this purpose, our *FacetCube* framework can be used as the underlying method for projecting documents into a given ontology space. Zaragoza et al. [36] proposed a Bayesian extension to the language model where prior knowledge is incorporated naturally as the hyper-parameter of the conjugate prior. Lafferty and Zhai [21] proposed a model for document retrieval using Bayesian decision theory, where prior knowledge and feedback can be incorporated through random walk on a Markov chain. Zhu et al. [39] incorporated

content and link information in a unified framework of matrix factorization for the task of document classification. In comparison, our work is an approach in the context of non-negative tensor factorization. Empirical comparison between these previous studies and our FacetCube framework in terms of effectiveness is one of our future directions.

2 Background

In this section, we introduce some background information as well as approaches based on the standard NTF.

2.1 Notations and definitions

First, we introduce some mathematical notations that will be used in this paper. In this paper, unless stated otherwise, values of all variables belong to \mathcal{R}_+ , the set of non-negative real numbers. We denote scalars by lower-case letters (e.g., a, b, α, β), vectors by lower-case letters in vector forms (e.g., \vec{p}, \vec{q}), matrices by capital letters (e.g., X, Y, Z), and tensors by calligraphic letters (e.g., $\mathcal{A}, \mathcal{B}, \mathcal{C}$). We reserve i, j, k, l, m, n to indicate the indices of tensors, and I, J, K, L, M, N for the sizes of the corresponding dimensions. In the following definitions and for the formulations in the rest of the paper, we restrict attention to 3rd-order tensors, whereas the same definitions and formulations can be easily generalized to tensors of arbitrary orders.

Dot product The dot product between tensors \mathcal{A} and \mathcal{B} is denoted by

$$\langle \mathcal{A}, \mathcal{B} \rangle = \sum_{ijk} \mathcal{A}_{ijk} \mathcal{B}_{ijk}$$

Similar dot products are defined on matrices and vectors. For example, $\langle U, V \rangle = \text{tr}(U^T V)$, where tr is the trace operator.

Element-wise product The element-wise product between tensors \mathcal{A} and \mathcal{B} is denoted by

$$\mathcal{C} = \mathcal{A} \circ \mathcal{B}$$

with $\mathcal{C}_{ijk} = \mathcal{A}_{ijk} \mathcal{B}_{ijk}$.

Element-wise division Similarly, the element-wise division is denoted by

$$\mathcal{C} = \mathcal{A} / \mathcal{B}$$

with $\mathcal{C}_{ijk} = \mathcal{A}_{ijk} / \mathcal{B}_{ijk}$. We define element-wise products and divisions on matrices and vectors in a similar way.

KL-divergence The KL-divergence between tensors \mathcal{A} and \mathcal{B} is defined as

$$KL(\mathcal{A} \| \mathcal{B}) = \sum_{ijk} \mathcal{A}_{ijk} \log \mathcal{A}_{ijk} / \mathcal{B}_{ijk} - \sum_{ijk} \mathcal{A}_{ijk} + \sum_{ijk} \mathcal{B}_{ijk}.$$

KL-divergence is defined on matrices and vectors similarly.

Base transform We say a tensor $\mathcal{A} \in \mathcal{R}_+^{I \times J \times K}$ is obtained from another tensor $\mathcal{C} \in \mathcal{R}_+^{L \times M \times N}$ by a base transform, if we have $\mathcal{A}_{ijk} = \sum_{lmn} X_{il} Y_{jm} Z_{kn} \mathcal{C}_{lmn}$. We denote the base transform by

$$\mathcal{A} = [\mathcal{C}, X, Y, Z].$$

Core tensor and facet matrices In the above base transform, we refer to \mathcal{C} as the core tensor and X , Y and Z as the facet matrices of the corresponding dimensions. The same notation applies to matrices where $[S, U, V] = USV^T$. In addition, we define shorthand notations for three special base transforms: $\mathcal{C} \times_1 X \doteq [\mathcal{C}, X, I_M, I_N]$, $\mathcal{C} \times_2 Y \doteq [\mathcal{C}, I_L, Y, I_N]$, and $\mathcal{C} \times_3 Z \doteq [\mathcal{C}, I_L, I_M, Z]$, where I_L, I_M, I_N are the identity matrices of sizes L, M , and N , respectively. (In the literature, \times_1, \times_2 and \times_3 are also called the mode-1, mode-2, and mode-3 multiplications of \mathcal{C} [9]).

Partial derivative of the inner product As the partial derivative of the inner product $\langle \mathcal{A}, [\mathcal{C}, X, Y, Z] \rangle$ with respect to X is independent to X , we denote it by $\langle \mathcal{A}, [\mathcal{C}, \bullet, Y, Z] \rangle$. That is, with $D_{il} = \sum_{jkmn} Y_{jm} Z_{kn} C_{lmn} A_{ijk}$, we have

$$\langle \mathcal{A}, [\mathcal{C}, \bullet, Y, Z] \rangle \doteq \frac{\partial}{\partial X} \langle \mathcal{A}, [\mathcal{C}, X, Y, Z] \rangle = D.$$

In a similar way, we use $\langle \mathcal{A}, [\mathcal{C}, X, \bullet, Z] \rangle$, $\langle \mathcal{A}, [\mathcal{C}, X, Y, \bullet] \rangle$, and $\langle \mathcal{A}, [\bullet, X, Y, Z] \rangle$ to indicate the partial derivatives of $\langle \mathcal{A}, [\mathcal{C}, X, Y, Z] \rangle$ with respect to Y , Z , and \mathcal{C} , respectively. Note that we have the following relation

$$\langle \mathcal{A}, [\bullet, X, Y, Z] \rangle \doteq \frac{\partial}{\partial \mathcal{C}} \langle \mathcal{A}, [\mathcal{C}, X, Y, Z] \rangle = [\mathcal{A}, X^T, Y^T, Z^T].$$

2.2 Data tensor, core tensor, and facet matrices

In an NTF-based approach that extracts data characteristics, the first step is to construct a *data tensor*, where the order of the data tensor is the same as the number of dimensions of the data. We again use the blogosphere as an example. A third-order data tensor $\mathcal{A} \in \mathcal{R}_+^{I \times J \times K}$ can be used to represent the blog data where the three dimensions of the data tensor correspond to blogger, keyword, and timestamp, respectively. Each of I, J , and K represents the size of the corresponding data dimensions. That is, there are in total I bloggers, J keywords and K timestamps in the data. Each entry of the data tensor represents the intensity of the corresponding entry in the observed data. For example, $(\mathcal{A})_{ijk} = a_{ijk}$ can be the frequency that blogger i mentioned keyword j at timestamp k .

Once the data tensor is constructed, in approaches based on the standard NTF, a non-negative tensor factorization is directly applied to the data tensor. The outcomes of this factorization consist of the two parts. The first part is a set of *facet matrices* where each matrix represents the most significant characteristics of one dimension of data. More specifically, the number of facet matrices is the same as the order of the tensor and for each facet matrix, each column of the matrix indicates one facet of the corresponding dimension of data. For our blog example, the facet matrices are $X \in \mathcal{R}_+^{I \times L}$, $Y \in \mathcal{R}_+^{J \times M}$, and $Z \in \mathcal{R}_+^{K \times N}$, where each column of X, Y , and Z denotes a salient community of bloggers, a significant topics in posts, and a noteworthy temporal trends, respectively.

The second part of the NTF decomposition is a *core tensor* \mathcal{C} , which has the same order as the data tensor but usually with a much smaller size. \mathcal{C} represents the correlation among all the facets in all the data dimensions. In our blog example, $\mathcal{C} \in \mathcal{R}_+^{L \times M \times N}$ where L, M , and N are the number of facets in the dimensions of blogger, keyword, and timestamp, respectively.

2.3 Approach based on the standard NTF

The target of a non-negative tensor factorization is to find the core tensor \mathcal{C} , the facet matrices X, Y , and Z , so that when put together as $[\mathcal{C}, X, Y, Z]$, they approximate \mathcal{A} in an optimal way. A commonly used metric to measure the approximation error is the KL-divergence.

That is, the objective of the standard NTF is to find C , X , Y , and Z to minimize the following error

$$\text{error}_{KL} = KL(\mathcal{A}||[C, X, Y, Z]). \quad (1)$$

In addition, because of the scale-free issue (e.g., $[C, X, Y, Z] = [\alpha C, \frac{1}{\alpha} X, Y, Z]$), additional constraints are added so that each column of X , Y , and Z must sum to 1, as well as the sum of all the entries of C must be 1.

The following procedure, which we proposed in our previous work [7, 8], can be used to iteratively searching the optimal solutions for C , X , Y , and Z .

Lemma 1 For a given $\mathcal{A} \in \mathcal{R}_+^{I \times J \times K}$, we first define an auxiliary tensor $\mathcal{B} \doteq \mathcal{A}/[C, X, Y, Z]$. Then the following update rules are guaranteed to converge to an optimal solution to the objective function defined in Eq. (1):

$$\begin{aligned} C &\leftarrow_0 C \circ \langle \mathcal{B}, [\bullet, X, Y, Z] \rangle, \\ X &\leftarrow_1 X \circ \langle \mathcal{B}, [C, \bullet, Y, Z] \rangle, \\ Y &\leftarrow_1 Y \circ \langle \mathcal{B}, [C, X, \bullet, Z] \rangle, \\ Z &\leftarrow_1 Z \circ \langle \mathcal{B}, [C, X, Y, \bullet] \rangle. \end{aligned}$$

where \leftarrow_0 denotes the operation of after all updates are completed, normalizing so that all entries sum to one, and \leftarrow_1 denotes the same operation except that all columns are normalized so that they sum to ones.

The detailed proof for Lemma 1 can be found in [8].

3 The FacetCube framework

In this section, we describe the details of our FacetCube framework. As we have mentioned in the introduction, incorporating users' prior knowledge can potentially benefit the extraction of data characteristics. The benefits can be of multiple-fold—the extracted characteristics can be more generalizable because of the alleviation of overfitting, they can be more reasonable because they fit users' prior knowledge, and they can meet the users' requirement when the users want to enforce the facets on certain data dimensions. To achieve these benefits, we extend the standard NTF by allowing users' prior knowledge to be incorporated in the process of decomposition. We will discuss two variants in our framework—when the facets in a data dimension are restricted to be in a user-given subspace and when the facets are fixed by the user. Together with the case of unconstrained facets in the standard NTF, our framework allows end users to control the process of data characteristic extraction at three different levels.

3.1 Basis-constrained data dimensions

In the first variant, we assume that the prior knowledge from the end users forms a subspace from which the facets can be located. To illustrate this variant, we first rewrite Eq. (1) as the following equivalent form

$$\text{error}_{KL} = KL(\mathcal{A}||[C, I_I X, I_J Y, I_K Z]), \quad (2)$$

where I_I , I_J , and I_K are identity matrices of dimensions I , J , and K , respectively. We can therefore consider, for example, every facet in Y as a point located in the simplex formed by

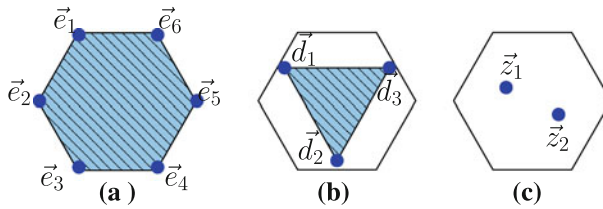


Fig. 1 Schematic illustrations of the FacetCube framework with **a** unconstrained, **b** basis-constrained, and **c** fixed data dimensions. In **a**, **b**, the shaded areas are the simplex in which the facets can locate; in **c**, \vec{z}_1 and \vec{z}_2 are fixed facets

$\{\vec{e}_1, \dots, \vec{e}_J\}$ where \vec{e}_j is the unit vector with dimension J (i.e., the j -th element being 1 and other elements being zeros). However, because J can be very large (i.e., tens of thousands of keywords), this simplex can have too large degree of freedoms, and as a result, the search space is too large for us to find reliable facets. In such a case, users' prior knowledge can be used to reduce the search space. Assuming the prior knowledge is given as a set of *basis* $Y_B = \{\vec{b}_1, \dots, \vec{b}_{M'}\}$, then we can reduce the search space of facets to be in the simplex formed by Y_B . That is, each facet is in the form of $Y_B \cdot \vec{y}$ for certain $\vec{y} \in \mathcal{R}^{M'}$.

To further clarify this, we come back to the blog example. Each facet in the content dimension is a significant topic and in the standard NTF approaches, a topic can be formed by any combination of keywords with arbitrary weights. If the end users provide a set of sub-topics, for example, obtained from the leaf nodes of a content ontology tree, then our framework can take these sub-topics as a *basis* for the content facets. That is, each facet (topic) must be a convex combination of the given sub-topics. In this case, we say the data dimension is *basis-constrained*. Figure 1b illustrates this point. Another example for basis-constrained data dimension can be the time dimension. Each facet in the time dimension corresponds to a noteworthy temporal trend. According to our prior intuition, a temporal trend should be smooth instead of noise-like. One way to incorporate this prior knowledge is to form a basis consisting of Fourier series in low frequencies (appropriately raised so that the values are non-negative). Then a convex combination of this basis will not contain high frequency components, and therefore, the facets will be guaranteed to be smooth.

In a mathematical form, to handle basis-constrained data dimensions, we extend the objective function in NTF to the following form

$$\text{error}_{KL} = KL(\mathcal{A}||[C, X, Y_B Y, Z]). \quad (3)$$

subject to each column of X , Y , and Z sums to 1. Note that in this objective function, without loss of generality, we assume that Y is the basis-constrained dimension that Y_B is the basis matrix provided by the end users, and that each column of Y_B sums to 1. As can be observed in Eq. (3), the term $Y_B Y$ replaces the term $I_J Y$ in Eq. (2) and so Eq. (2) can be considered as a special case of Eq. (3).

3.2 Fixed data dimensions

In some other cases, the end users may require the facets for certain data dimensions to be *fixed*. Such a requirement may sound overly restricted, but it is not uncommon in real applications. For example, from the top-level nodes of an ontology tree, the user may have already determined the *top-level topics* (in contrast to the *sub-topics* in Sect. 3.1) such as politics, technologies, sports, etc., as well as the representation of these topics as facets. And the user's goal is to summarize the blog posts through the lens of these pre-constructed facets.

As another example, in the time dimension, the user may choose to use a set of facets that correspond to the domain knowledge, that is, concepts such as year, quarter, month, etc., to detect seasonal trends. In such scenarios, we extend the objective function in NTF to the following form

$$\text{error}_{KL} = KL(\mathcal{A}||[C, X, Y, Z_B]). \quad (4)$$

subject to each column of X , and Y sums to 1. Note that without loss of generality, we assume that the facets in the time dimension are fixed as Z_B and each column of Z_B sums to 1. Because Z_B is fixed, we say that Z is a *fixed data dimension*. Figure 1c illustrates this case.

3.3 The overall picture

To summarize the above two variants as well as the third variant which corresponds to the traditional NTF case, we give the objective function in its most general form as the following

$$\text{error}_{KL} = KL(\mathcal{A}||[C, X_B X, Y_B Y, Z_B Z]), \quad (5)$$

where X_B , Y_B , and Z_B are given a priori whereas X , Y , Z , and C are to be computed. When X_B , Y_B , and Z_B are set to identity matrices, we obtain the standard NTF problem; when any of the X_B , Y_B , or Z_B is given and the corresponding X , Y , or Z is to be computed, we have the case of the basis-constrained data dimension; when any of the X_B , Y_B , or Z_B is given and the corresponding X , Y , or Z is fixed to be an identity matrix, we have the case of the fixed data dimension.

In addition, because in our proposed algorithm, as will be described momentarily, the (per iteration) optimization of each data dimension is performed independently of other data dimensions, we can adopt any of the above three variants in our FacetCube framework for any data dimensions independently. For example, when using our framework, the end users may decide to (1) let the blogger communities be discovered with no constraints, (2) provide a basis for the contents, and (3) fix the facets in the time dimension.

3.4 An iterative algorithm

We propose an algorithm for computing optimal solutions for our FacetCube framework. We present the algorithm in the most general form—the basis-constrained form—and then show that the unconstrained and the fixed dimensions can be solved by using the same result.

The iterative algorithm and its correctness are stated in the following theorem.

Theorem 1 For a given $\mathcal{A} \in \mathcal{R}_+^{I \times J \times K}$, the following update rules are guaranteed to converge to an optimal solution to the objective function defined in Eq. (5).

$$B \leftarrow \mathcal{A} / [C, X_B X, Y_B Y, Z_B Z], \quad (6)$$

$$C \leftarrow_0 C \circ \langle B, [\bullet, X_B X, Y_B Y, Z_B Z] \rangle, \quad (7)$$

$$X \leftarrow_1 X \circ \langle B \times_1 X_B^T, [C, \bullet, Y_B Y, Z_B Z] \rangle, \quad (8)$$

$$Y \leftarrow_1 Y \circ \langle B \times_2 Y_B^T, [C, X_B X, \bullet, Z_B Z] \rangle, \quad (9)$$

$$Z \leftarrow_1 Z \circ \langle B \times_3 Z_B^T, [C, X_B X, Y_B Y, \bullet] \rangle. \quad (10)$$

where \leftarrow_0 and \leftarrow_1 are the same as those defined in Lemma 1.

The detailed proof for Theorem 1 is provided in the “Appendix”.

Although the iterative algorithm is given for the case of the basis-constrained dimension, it is straightforward to adopt it to the other two cases. For the case of unconstrained dimension, say $Y_B = I_J$ is the identity matrix, we directly use \mathcal{B} instead of $\mathcal{B} \times_2 Y_B^T$ in Eq. (9) of the update algorithm. For the case of fixed dimension, say $Z_B Z = Z_B$, we simply skip the corresponding update step in Eq. (10). Such a skip is safe because in each step of the iterative algorithm, Eqs. (7)–(10) *independently* improve the objective function in a monotonic way.

It is worth noting that another consequence of such independent updates is that the end users can simultaneously enforce different levels of preference at each dimension, for example, by setting $X_B = I_I$ (unconstrained), providing Y_B (basis-constrained), and fixing $Z_B Z = Z_B$ (fixed) at the same time.

4 A probabilistic interpretation

The FacetCube framework turns out to have a natural probabilistic interpretation in that it is equivalent to a special probabilistic generative model. In the following discussion, we focus on the basis-constrained dimension, and we use the blog data as an example, assuming the observed data are in a list in the form of $\{(\text{blogger } i, \text{word } j, \text{time } k, \mathcal{A}_{ijk})\}$. We can use the following probabilistic generative procedure to describe how the observed data are sampled:

1. select a community l , a topic m , a temporal trend n , with probability c_{lmn} (this prior probability corresponds to an entry of the core tensor \mathcal{C}),
2. conditioning on the result in step 1, select a sub-community l' , a sub-topic m' , a sub-trend n' , following $p(l'|l)$, $p(m'|m)$, and $p(n'|n)$ (these conditional probabilities correspond to the l -th, m -th, and n -th columns of X , Y , and Z , respectively),
3. conditioning on the results in step 2, select a blogger i , a word j , a time stamp k , following $p(i|l')$, $p(j|m')$, and $p(k|n')$ (these conditional probabilities correspond to the l' -th, m' -th, and n' -th columns of X_B , Y_B , and Z_B , respectively).

Then under this generative model, the log-likelihood of the data can be written as

$$\sum_{ijk} \mathcal{A}_{ijk} \log \left(\sum c_{lmn} p(l'|l) p(m'|m) p(n'|n) p(i|l') p(j|m') p(k|n') \right) \quad (11)$$

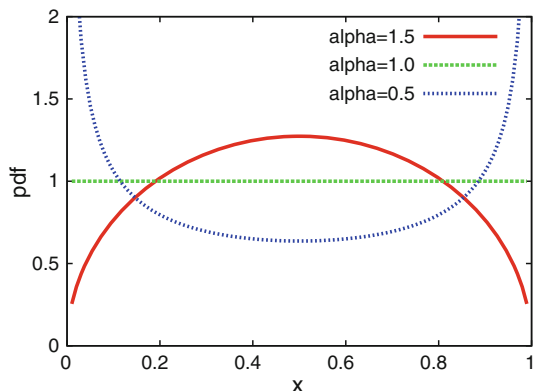
where the inner sum is computed over all l, m, n, l', m', n' . A simple derivation can show that maximizing the log-likelihood in Eq. (11) is equivalent to minimizing the KL loss in Eq. (5). As a consequence, our FacetCube framework is equivalent to this probabilistic generative model.

The probabilistic interpretation, other than giving additional insights to our FacetCube framework, provides the underpinning for certain extensions of our basic framework. In the remainder of this section, we show such an extension which is to consider a maximum a posterior (MAP) estimation for the basic probabilistic model. Here, we consider a case of using the Dirichlet distribution as the prior distribution of the parameters. First, we consider the prior for X . The prior of each column of X is a Dirichlet distribution with hyper-parameter $\alpha_X > 0$. The logarithm of the prior probability is

$$\ln P(X) = (\alpha_X - 1) \sum_{il} \ln X_{il} + c_X,$$

where c_X is a value irrelevant to X . Similarly, we assume the priors for Y, Z and \mathcal{C} are all Dirichlet distributions, with hyper-parameters α_Y, α_Z , and α_C , respectively.

Fig. 2 Distribution of variables with different α . When $\alpha > 1$, the mode is at 0.5, and so x tends to be non-zeros. When $\alpha < 1$, the modes are 0 and 1, which make x or $1 - x$ tend to be 0, and so it results in sparse solutions



The logarithm of the error for the MAP estimation is that for the MLE plus the logarithm of the prior probabilities:

$$\text{error}_{\text{MAP}} \doteq \text{error}_{KL} - \ln P(X) - \ln P(Y) - \ln P(Z) - \ln P(C). \quad (12)$$

With a few derivations, we can obtain the following algorithm for solving the MAP estimation.

Corollary 1 *The following update rules are guaranteed to converge to an optimal solution to the objective function defined in Eq. (12).*

$$\begin{aligned} \mathcal{B} &\leftarrow \mathcal{A}/[\mathcal{C}, X_B X, Y_B Y, Z_B Z], \\ \mathcal{C} &\leftarrow_{0,\epsilon} \mathcal{C} \circ \langle \mathcal{B}, [\bullet, X_B X, Y_B Y, Z_B Z] \rangle + (\alpha_{C-1}) \\ X &\leftarrow_{1,\epsilon} X \circ \langle \mathcal{B} \times_1 X_B^T, [\mathcal{C}, \bullet, Y_B Y, Z_B Z] \rangle + (\alpha_X - 1), \\ Y &\leftarrow_{1,\epsilon} Y \circ \langle \mathcal{B} \times_2 Y_B^T, [\mathcal{C}, X_B X, \bullet, Z_B Z] \rangle + (\alpha_Y - 1), \\ Z &\leftarrow_{1,\epsilon} Z \circ \langle \mathcal{B} \times_3 Z_B^T, [\mathcal{C}, X_B X, Y_B Y, \bullet] \rangle + (\alpha_Z - 1). \end{aligned}$$

where ϵ is a small positive real number and $\leftarrow_{0,\epsilon}$ denotes the operation of after all updates are completed, applying projection \mathcal{P}_ϵ on all entries as a vector, and $\leftarrow_{1,\epsilon}$ denotes the operation of after all updates are completed, applying projection \mathcal{P}_ϵ on each column, and projection \mathcal{P}_ϵ is

$$\mathcal{P}_\epsilon^p \zeta = \arg \min_{\xi: \xi_k \geq \epsilon \text{ and } \sum_k \xi_k = 1} - \sum_k \zeta_k \ln \xi_k.$$

The detailed proof for Corollary 1 and efficient algorithm for \mathcal{P}_ϵ are provided in the “Appendix”.

Corollary 1 reveals that the Dirichlet prior can play two different roles in computing the optimal solutions. When $\alpha > 1$, the Dirichlet prior plays a role of *smoothing* by adding pseudo-counts to each variable. On the other hand, when $\alpha < 1$, the Dirichlet prior plays a role opposite to smoothing—it pulls variable with small values into the negative territory (which are then set to ϵ by the update rules) and as a result, making the solution more *sparse* (if we disregard the small-valued ϵ 's). This point is demonstrated in Fig. 2 by using a two-variable Dirichlet (Beta) distribution.

5 Implementation of algorithm

In this section, we first describe several techniques we developed to make our algorithm scalable and then introduce a technique we adopted for fast computation of top- K queries. The reason for us to give such a detailed description is that we believe these techniques are important by themselves: they can be equally applied to other multiplicative update rules for tensors of arbitrary order and they can be used to solve a family of similar problems.

5.1 Scalable implementation

A naive implementation of the update rules in Theorem 1 is impractical for all but very small data sets. In real applications, the size of a data dimension (e.g., the number of bloggers or the size of the vocabulary) can easily be tens of thousands. When $I = J = K = 10,000$, for example, even the computation of $[\mathcal{C}, X, Y, Z]$ in Eq. (13) is unmanageable, for $[\mathcal{C}, X, Y, Z]$ is a dense tensor with 1 trillion entries.² To handle data from real applications, we designed a scalable implementation of the FacetCube algorithm that exploits three key insights we observed. Part of these insights has been roughly mentioned in a high-level way in [7,8]. Here we give them a much more rigorous and complete treatment.

In the following discussion, without loss of generality, it is assumed that $L \geq M \geq N$ and $I \geq J \geq K$. In addition, we assume that each data record is stored in the format of $(key1, key2, key3, v)$, where $key1, key2$ and $key3$ are the indices of the data record in the first, second, and third dimensions, respectively; v is the value of the data record, which can be, for example, the number of times that a blogger mentions a keyword on a particular day. It is also assumed that the data records have been sorted according to $key3$ as the major key and then $key2$ and then $key1$ as the minor keys. We shall come back to this last assumption later.

5.1.1 Taking advantage of data sparseness

The first insight we have exploited is that data in real applications are usually sparse. For example, not every blogger uses every word in every of his or her posts. Because \mathcal{A} is a sparse tensor, not every entry of $[\mathcal{C}, X, Y, Z]$ is needed. In our implementation, we take advantage of data sparseness by computing base transforms in an *on-demand* fashion.

More specifically, because the update rules in Theorem 1 mainly involve *element-wise* multiplications and divisions and because \mathcal{A} and \mathcal{B} are both sparse, in our algorithm we never explicitly compute the entries in $[\mathcal{C}, X, Y, Z]$ that do not correspond to non-zero values in \mathcal{A} and \mathcal{B} . It can be shown that it only takes time $O(n_z \cdot L^3)$ to compute the entries in $[\mathcal{C}, X, Y, Z]$ that correspond to the non-zero entries in \mathcal{A} , where n_z is the total number of non-zero entries in \mathcal{A} .

5.1.2 Caching intermediate results

The second insight is that the update rules in Theorem 1 involve a lot of nested summations and multiplications that have a lot of structures. Therefore if we carefully order these nested operations and cache the intermediate results, we can avoid generating them multiple times. More specifically, we rewrite Eq. (13) as the following by pushing certain summations inside the nest

² For the discussion in this part, we restrict attention to the case of unconstrained dimensions, because the computations for $X_B X$, $Y_B Y$, and $Z_B Z$ do not dominate the time complexity.

$$\mathcal{B}_{ijk} = \mathcal{A}_{ijk} / \sum_l X_{il} \left(\sum_m Y_{jm} \left(\sum_n Z_{kn} C_{lmn} \right) \right). \quad (13)$$

In the above expression, for entries with the same k index, the term in the inner parentheses can be reused even when the i and j indices vary; similarly, for entries with the same j and k indices, the term in the outer parentheses can be reused even when the i index varies.

We can derive similar expressions for Eqs. (7)–(10). For instance, Eq. (7) can be rewritten as

$$[\mathcal{B}, X^T, Y^T, Z^T]_{lmn} = \sum_k Z_{kn} \left(\sum_j Y_{jm} \left(\sum_i X_{il} \mathcal{B}_{ijk} \right) \right). \quad (14)$$

And for Eq. (8), we first notice that

$$\langle \mathcal{B} \times_1 X_B^T, [\mathcal{C}, \bullet, Y, Z] \rangle = X_B^T \langle \mathcal{B}, [\mathcal{C}, \bullet, Y, Z] \rangle$$

and for the second part of the above expression, we have

$$\langle \mathcal{B}, [\mathcal{C}, \bullet, Y, Z] \rangle_{il} = \sum_{jk} \mathcal{B}_{ijk} \left(\sum_m Y_{jm} \left(\sum_n Z_{kn} C_{lmn} \right) \right) \quad (15)$$

As a result, we can reuse intermediate computation results in all the update rules, since \mathcal{A} is stored with *key3* (the k index) as the major key and then *key2* (the j index) and then *key1* (the i index) as the minor keys. The sorting of \mathcal{A} does affect the time complexity for two reasons. On the one hand, the indices i , j , and k are positive integers with known upper-bounds and so a linear sorting algorithm, such as bucket sort, can be applied. On the other hand, \mathcal{A} only has to be sorted once before the iterative process starts and so the sorting cost is amortized among multiple iterations. It is worth mentioning that \mathcal{B} , which is also sparse, does not have to be explicitly sorted—it is automatically sorted because of the way it is derived from the sparse tensor \mathcal{A} .

5.1.3 Ordering the dimensions

The third insight is that different data dimensions usually have different cardinalities. For example, if the time unit is a day, then the size of the time dimension can be much smaller than that of the blogger dimension. We take advantage of this fact by ordering the computation in such a way that the dimension of the smallest size is put in the inner parentheses. A moment of thoughts on Eqs. (13)–(15) can show that this ordering is always possible as long as \mathcal{A} is sorted accordingly at the beginning. As will be demonstrated in the experimental studies, this re-ordering makes huge difference in the running time of our algorithm.

5.1.4 Putting it all together

To summarize, as an example, Fig. 3 gives the high-level pseudo code for our implementation of the update rule for updating X , for the case of basis-constrained dimension. The time complexity can be shown to be $O(n_z \cdot L + n_p \cdot L^2 + K \cdot L^3)$, where n_z is the number of non-zero entries in \mathcal{A} , n_p is the number of distinct (j, k) pairs in \mathcal{A} , and K is the smallest size among all the data dimensions. As can be observed from the time complexity, if we consider the number of factors L as a constant, then the per iteration time is linear in the size of raw data n_z ; if on the other hand we do not consider L to be a constant, because $n_z > n_p > K$,

Algorithm for updating X	
▷	input: \mathcal{B} as $\{\langle key1, key2, key3, v \rangle\}, X, Y, Z, \mathcal{C}, X_B$
▷	output: updated X
1:	$k \leftarrow -1, j \leftarrow -1, i \leftarrow -1, E \leftarrow \mathbf{0};$
2:	for each entry $\langle key1, key2, key3, v \rangle$ of \mathcal{B} do
3:	if $k \neq key3$
4:	$k \leftarrow key3, j \leftarrow -1;$
5:	construct D s.t. $D_{lm} \leftarrow \sum_n Z_{kn} \mathcal{C}_{lmn};$
6:	if $j \neq key2$
7:	$j \leftarrow key2;$
8:	construct \vec{d} s.t. $(\vec{d})_l \leftarrow \sum_m Y_{jm} D_{lm};$
9:	$i \leftarrow key1;$
10:	$E_{row_i} \leftarrow E_{row_i} + v \cdot \vec{d}^T;$
11:	$X \leftarrow_1 X \circ (X_B^T E);$
12:	return $X;$

Fig. 3 Implementation of the update rule for X . Those for Y , Z , and \mathcal{C} are similar

we assign the lowest coefficient to the most expensive L^3 term and therefore fully exploit the data characteristics. We will verify these analyses in the experimental studies.

5.2 Fast computation for top- K queries

In many recommendation applications, instead of being interested in the ranks of the whole list of candidates, the end users often are only interested in a quick answer to the top- K queries. For example, the query may be “Who are the top 3 bloggers mostly involved in a given topic during a given time period?” or “What are the top 10 references that are mostly relevant to a given group of authors who plan to co-author a paper on a given set of keywords?”. Because in many applications, such queries must be answered in real time, fast computation of top- K answers becomes crucial. In the implementation of our framework, we develop a technique that derives the *exact* top- K answers but usually does not compute the scores of *all* the candidates for the recommendation.

5.2.1 Condition of monotonicity

The key component of the technique is to use the “threshold algorithm” (TA) developed by [13] which has been extensively used in the database field for answering fuzzy queries over multimedia databases. For the TA algorithm to work, a key requirement is that the score function must be *monotonic*. A function $f(z_1, \dots, z_N)$ is monotonic if $\{z'_1 \geq z_1, \dots, z'_N \geq z_N\}$ implies that $f(z'_1, \dots, z'_N) \geq f(z_1, \dots, z_N)$. It turns out that our ranking function satisfies this condition of monotonicity.

We illustrate this by using the above query of recommending references to a given set of authors on a given set of keywords. We assume that X, Y , and Z correspond to *author*, *keyword*, and *reference*, and X_B, Y_B , and Z_B are the corresponding basis. Recall that $\mathcal{A} \sim [C, X_B X, Y_B Y, Z_B Z]$ and so for a set of authors and a set of keywords, the relevance of the references are $[C, \vec{x}^T, \vec{y}^T, Z_B Z] = (Z_B Z) \cdot [C, \vec{x}^T, \vec{y}^T, I_N]$, where the k -th row of the result indicates the relevance of the k -th reference. Because \vec{x}^T and \vec{y}^T are obtained by aggregating the given set of authors and the given set of keywords, respectively, entries in

$\vec{c} = [C, \vec{x}^T, \vec{y}^T, I_N]$ are non-negative. Because $Z_B Z$ is also non-negative, it can be easily shown that the score function $f[(Z_B Z)_{row_k}] = (Z_B Z)_{row_k} \cdot \vec{c}$ is a monotonic function.

5.2.2 Efficient top- K recommendations

In this subsection, we describe how the top- K recommendations are computed in our implementation. However, because we have already use K to represent the size of the dimension of Z , in this subsection, we assume the top- r recommendations are to be made.

We first describe some data structures used in our implementation. A *pair* p is a real-integer pair $(p.value, p.rid)$. A priority queue Q is a priority queue of *pairs*, with the priority determined by $p.value$. We assume Q supports normal operations on a priority queue, including $Q.size()$, $Q.top()$, $Q.pop()$, and $Q.insert(p)$. Note that $Q.top()$ is the *pair* p in Q that has the minimal $p.value$. The last data structure P is a $K \times N$ matrix of *pairs*, where the (i, j) entry of P represents the i -th largest value in the j -th column of $Z_B Z$, together with the corresponding row id. More specifically, we first get the j -th column of $Z_B Z$, sort it by the entry values, then put the sorted values (together with the original row ids, to form *pairs*) into the j -th column of P . It is worth pointing out that P only need to be computed once, probably in an offline fashion, and then can be used for all the recommendation computations.

Figure 4 provides the details on how the top- r recommendations are efficiently computed in our implementation. The key point is that line 26 can be reached much earlier before the scores for *all* the references have been computed. We will illustrate the effectiveness of this algorithm in experimental studies.

6 Experimental studies

In this section, we present experimental studies on two data sets—a paper citation data set and a blog data set. For experiments on the paper citation data, we study the task of personalized reference recommendation, and we focus on the accuracy of recommendations, the running time of our factorization algorithm, the time for computing the top- K recommendations, as well as the impact of the Dirichlet prior. For experiments on the blog data, we illustrate a real application scenario where the end users have an ontology as the facets for content, and a fixed set of facets for time.

6.1 The CiteSeer data

The paper citation data set is obtained from the CiteSeer website.³ This data set contains the information about papers published in the computer science area. For a detailed description of the data set, we refer interested readers to [38]. The information we used includes the authors, the abstract, and the citations of each paper. The abstracts of the papers are pre-processed through the following steps. First, an NLP package⁴ is used to detect and only keep the nouns from the abstracts (according to our experiences, non-noun words are not very informative); then a standard set of stop words are removed; finally, the Porter stemming algorithm is applied to normalize the keywords. In all these steps, out-of-box models are used directly without trained by other data. After these steps, we use certain thresholds to filter out

³ <http://citeseer.ist.psu.edu/>.

⁴ <http://opennlp.sourceforge.net/>.

Algorithm for making top- r recommendations

```

▷ input:  $Z_B Z, P, \vec{c}, r$ 
▷ output:  $Q$ , the set of top- $r$  recommendations
1:  $t \leftarrow 0.0$ ;
2: for  $j \leftarrow 1$  to  $N$  do
3:    $t \leftarrow t + \vec{c}_{row_j} \cdot P(1, j).value$ ;
4:    $id \leftarrow P(1, j).rid$ ;
5:   if  $id$  has not been seen before
6:      $v \leftarrow (Z_B Z)_{row_{id}} \cdot \vec{c}$ ;
7:     if  $Q.size() < r$ 
8:        $Q.insert(pair(id, v))$ ;
9:     else if  $v > Q.top().value$ 
10:       $Q.pop()$ ;
11:       $Q.insert(pair(id, v))$ ;
12:  $i \leftarrow 1$ ;
13: while  $i < K$  and ( $Q.size() < r$  or  $Q.top().value < t$ ) do
14:    $i \leftarrow i + 1$ ;
15:   for  $j \leftarrow 1$  to  $N$  do
16:      $t \leftarrow t + \vec{c}_{row_j} \cdot (P(i, j).value - P(i - 1, j).value)$ ;
17:      $id \leftarrow P(i, j).rid$ ;
18:     if  $id$  has not been seen before
19:        $v \leftarrow (Z_B Z)_{row_{id}} \cdot \vec{c}$ ;
20:       if  $Q.size() < r$ 
21:          $Q.insert(pair(id, v))$ ;
22:       else if  $v > Q.top().value$ 
23:          $Q.pop()$ ;
24:          $Q.insert(pair(id, v))$ ;
25:       if  $Q.size() = r$  and  $Q.top().value \geq t$ 
26:         return  $Q$ ;
27: return  $Q$ ;

```

Fig. 4 Computing top- r recommendations in an efficient way by using the TA algorithm

uncommon words, authors with few publications, and references being cited too infrequently. The final data set contains 15,072 papers with 6,821 distinct authors, 1,790 distinct keywords, and 21,894 distinct references.

We design the task as for a set of authors and a set of keywords, to recommend the most relevant references. Such recommendations can be helpful for authors to select a set of candidate references when writing a paper on a particular topic. The performance is measured in the following way. The 15,072 papers are randomly split evenly into a training set and a testing set with a 50–50 ration. Papers in the training set are used to construct the data tensor where the three data dimensions are: *author*(X), *keyword*(Y), and *reference*(Z). The standard NTF and FacetCube are applied to the data tensor to compute the facet matrices and the core tensor. During testing, for each paper in the testing set, we assume its authors and keywords are given. Relevant references are to be recommended by using $[C, \vec{x}^T, \vec{y}^T, Z] = Z \cdot [C, \vec{x}^T, \vec{y}^T, I_N]$, where \vec{x}^T and \vec{y}^T are obtained by aggregating the authors and keywords, respectively, of the given paper. The performance is measured by comparing the recommendations with the ground truth, namely the real references cited in the papers in the testing set.

Table 1 Performance of different algorithms on the task of recommending references on CiteSeer data

	Top-1	Top-5	Top-10	Top-50	Top-100
Baseline	0.113	0.316	0.438	0.810	1.007
NTF	0.109	0.322	0.456	0.864	1.081
FacetCube	0.110	0.326	0.463	0.870	1.088
FacetCube-D	0.110	0.325	0.462	0.868	1.086

The metric we used for the measurement is the average Normalized Discounted Cumulative Gain (NDCG) [19] for the top- K answers, where NDCG is defined as

$$\text{NDCG} = \sum_{k=1}^K \frac{2^{\text{rel}(k)} - 1}{\log(1 + k)}.$$

NDCG is a relatively new measure used extensively in Web search. It measures the quality of a ranked list (obtained by a ranking algorithm) by the usefulness, or gain, of documents based on their positions in the ranked list. It allows different levels of relevance and weighs the recommendations according to their ranks in the ranked list. In our performance study, we choose to use a binary relevance score for $\text{rel}(k)$, that is, $\text{rel}(k) = 1$ if item k occurs in the testing data and 0 otherwise.

For the FacetCube algorithm, we show a case where the prior knowledge is used to improve the personalized recommendation. We construct a basis X_B for the authors by using the co-authorship relationship. More specially, we set X_B to be the normalized version of $I + \alpha W$ where I is the identity matrix and W is the adjacency matrix representing the co-authorship relations in the training data, that is, $W_{ij} = 1$ if authors i and j have co-authored a paper and 0 otherwise. X_B can actually be considered as the first-order approximation to the diffusion kernel $(I - \alpha W)^{-1}$ and therefore by using X_B , we restrict the facets of authors to be embedded in the manifold of this diffusion kernel. In our experiments, we simply set $\alpha = 0.01$.

We compare FacetCube with (a) the standard NTF and (b) a baseline in which the counts of reference for each keyword in the training data are stored and used to rank the references in the testing data. The baseline gives global (non-personalized) recommendations. For both the standard NTF and FacetCube, we set the facet numbers to be $L = M = N = 50$. The reason for us to choose 50 is that this relatively large data set contains a large number of facets (i.e., diversified topics and large number of clusters of authors); we set $L = M = N$ rather arbitrarily because we do not have other insights. Table 1 reports the NDCG scores at top- K , for $K = 1, 5, 10, 50, 100$. As can be seen from the results, FacetCube slightly outperforms the standard NTF in all the cases and outperforms the baseline in most of the cases. FacetCube improves over the standard NTF because intuitively, we enforce the authors who have co-authored papers in the training set to be in similar facets in the factorization. However, the improvement of FacetCube over NTF is not very significant, which suggests that co-authorship is not a very strong indicator in terms of paper citation. In other words, the papers that an author tends to cite depend mainly on the author's background (e.g., which papers he or she is familiar with) and less on who the co-authors are.

To investigate the running time of our implementation, we conduct the following experiments. First, we construct a series of data sets with different sizes by randomly removing a portion of references. Figure 5a shows the running time of FacetCube under different data sizes. For the factorizations, we set the condition for convergence as the relative error being less than 10^{-4} , for which all the methods converge within a couple hundreds of

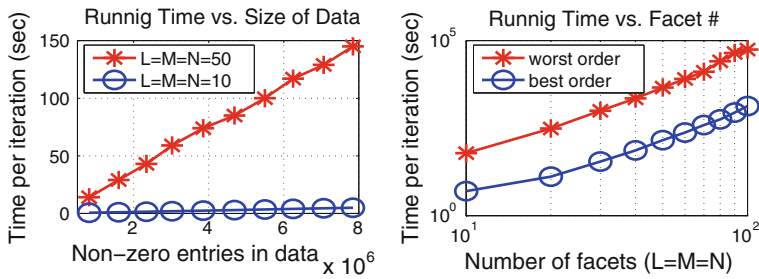


Fig. 5 Running time versus number of non-zero entries in the data tensor (a), and versus the number of facets (b)

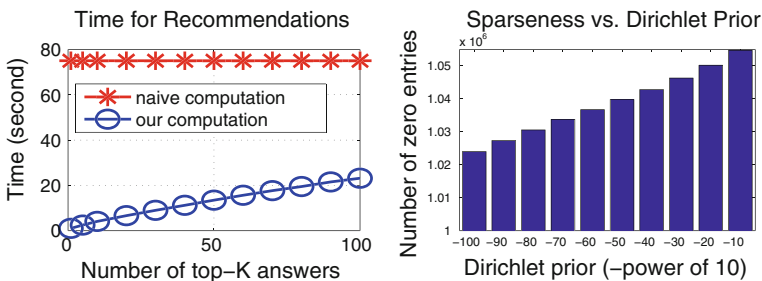


Fig. 6 Time for computing the top- K answers (a) and the level of sparsity (b)

iterations. The results in Fig. 5a verify that the running time of our implementation is linear in the size of data set. Second, we fix the data set and test a wide range of facet numbers, from 10 to 100. In Fig. 5b we show the running time under different facet numbers for two ways of ordering data: one by $\langle \text{reference}, \text{author}, \text{keyword} \rangle$ (best) and another by $\langle \text{keyword}, \text{author}, \text{reference} \rangle$ (worst). The former is the best because in this data set, the number of distinct keywords is much smaller than that of the references. From the curves, we can see that the running time is about $O(L^3)$ (note the log-log scale in the figure). Also revealed is that a good order can make the implementation 10 to 50 times faster for this specific data set.

In Fig. 6a, we report the time for computing the top- K queries by our implementation. Also shown is the lower bound of a naive implementation that computes the scores for all the candidates (for the naive approach, we only measure the time for computing the scores but not for ranking, and so it is a lower bound). As can be seen, our computation is much faster. For example, for the top-1 query, a naive approach will take at least 75 seconds for the testing data while our implementation takes less than 1 second. Next, we study how the Dirichlet prior influences the sparseness of the results. We run FacetCube with $\varepsilon = 10^{-100}$ and $\alpha_Z - 1$ to have a wide range of values. We report in Fig. 6b the number of entries in the facet matrix Z that are smaller than or equal to ε . As can be seen, with smaller values of $\alpha_Z - 1$, we make the facets more sparse. However, this sparseness has its costs. In Table 1, on the row of FacetCube-D, we give the NDCG scores for FacetCube with $\alpha_Z - 1$ set to be -10^{-50} . As can be seen, the performance deteriorates a little bit.

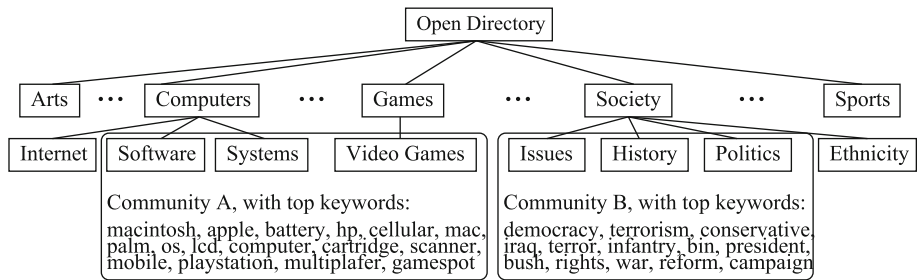


Fig. 7 Two representative communities and their projections on the Open Directory ontology

6.2 The blog data

The blog data were collected by an NEC in-house blog crawler. At the NEC Laboratories American, we have built a blog crawler with technology focuses. Due to space limit, we only give a high-level description of the crawler. There are two databases used by the crawler. The first database contains a set of “seed blogs”, which initially consist of some well-known blogs with technology focuses. For the seed blogs, the crawler continuously crawls their RSS feeds and then the corresponding entries. For each newly crawled entry, its content is analyzed and the hyperlinks embedded in the content are extracted. If an extracted hyperlink points to another entry and that entry belongs to a blog who is not a member of the seed blogs, then that entry and its blog are stored into the second database. The second database is checked regularly to see whether any blog in the database meets the criteria to become a new seed blog (the criteria are based on the number of citations and trackbacks from current seed blogs) and if so, that blog is moved to the first database and starts to be crawled continuously. This data set contains 1,161 bloggers, 280,259 blog posts, with 19,390 keywords (nouns only), for 28 consecutive weeks, between March 02 to September 09, 2008. In addition, we constructed an ontology for content by using data from the Open Directory Project (<http://www.dmoz.org/>). The ontology consists of 296 tree nodes, such as “Computers_Hardware” and “Society_Politics”. For each of the tree node, a description vector is constructed by aggregating the descriptions of all the Web sites under the given tree node.

In this experiment we demonstrate a scenario where the user’s domain knowledge is used to directly control the factorization in the FacetCube framework. For this purpose, we fix the facets in the content dimension to be the 296 description vectors and the facets in the time dimension to be 28 windowing functions that correspond to the 28 weeks in our data. Therefore, the free variables in this experiment are the blogger facet matrix X and the core tensor \mathcal{C} . We set the number of blogger facets (i.e., the number of blogger communities) to be 10.

It turns out that the blogger communities derived by the FacetCube algorithm all have certain focused topics, such as technology, politics, music, etc. In Fig. 7, we illustrate two representative blogger communities and their projections on the Open Directory ontology. The projections are computed by aggregating out the time dimension of the core tensor \mathcal{C} and from the resulting matrix, identifying for each community the top 3 ontology tree nodes that have the largest values. As can be seen, community *A* focuses on different issues related to technology while community *B* focuses on political issues (as verified by the top keywords among the top bloggers in the communities).

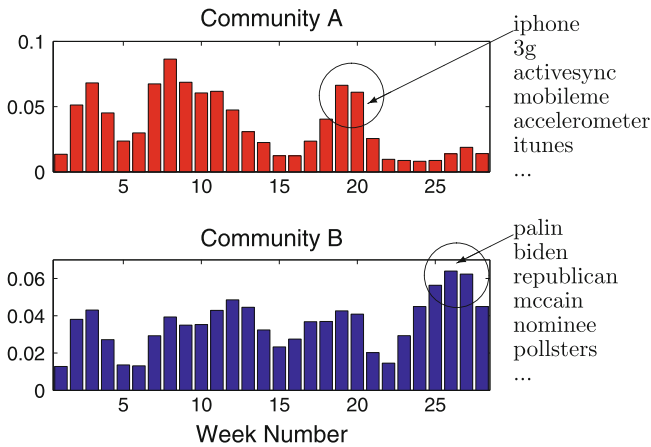


Fig. 8 Temporal trends for the two representative communities

Finally in Fig. 8, we show the temporal trends for these two communities by aggregating out the content dimension from the core tensor \mathcal{C} . In addition, we also show some hot keywords among the community members during certain temporal peaks. As can be seen, these hot keywords clearly indicate certain noteworthy events (e.g., the release of 3G iPhone in community A and the nomination of the vice presidential candidates in community B).

7 Conclusions

In this paper, we presented a novel framework, FacetCube, that extends the standard non-negative tensor factorization for extracting data characteristics from polyadic data. The FacetCube framework allows end users great flexibility in incorporating their prior knowledge at different levels in the process of characteristic extraction. In addition to describing the new framework from the NTF point of view, we also made connection to the perspective of probabilistic generative models. Furthermore, we provided iterative algorithms for solving the corresponding optimization problems as well as efficient implementations for handling large-scale data from real-life applications. Extensive experimental studies demonstrated the effectiveness of our new framework as well as the efficiency of our algorithm.

Acknowledgments The authors would like to thank Professor C. Lee Giles for providing the CiteSeer data set and thank Koji Hino and Junichi Tatemura for helping prepare the blog data set.

Appendix

Proof for Theorem 1

Proof Assume that the values obtained from the previous iteration are \tilde{X} , \tilde{Y} , \tilde{Z} , and $\tilde{\mathcal{C}}$, respectively. We prove the update rule for X . The rules for Y , Z , and \mathcal{C} can be proved similarly. For the update rule of X , we can consider Y , Z , and \mathcal{C} as fixed (i.e., fixed as their values \tilde{Y} , \tilde{Z} , and $\tilde{\mathcal{C}}$ in the previous iteration). To avoid notation clutters, we define $\tilde{\tilde{X}} \doteq X_B \tilde{X}$, $\tilde{\tilde{Y}} \doteq Y_B \tilde{Y}$, $\tilde{\tilde{Z}} \doteq Z_B \tilde{Z}$, and we rewrite the objective function as

$$\min D_X(X) = \min KL(\mathcal{A} || [\tilde{C}, X_B X, \tilde{Y}, \tilde{Z}])$$

First define

$$\gamma_{ijklmnl'} = \tilde{C}_{lmn}(X_B)_{il'} \tilde{X}_{l'l} \tilde{Y}_{jm} \tilde{Z}_{kn}$$

and

$$\theta_{ijklmnl'} = \frac{\tilde{C}_{lmn}(X_B)_{il'} \tilde{X}_{l'l} \tilde{Y}_{jm} \tilde{Z}_{kn}}{[\tilde{C}, \tilde{X}, \tilde{Y}, \tilde{Z}]_{ijk}} = \frac{\gamma_{ijklmnl'}}{[\tilde{C}, \tilde{X}, \tilde{Y}, \tilde{Z}]_{ijk}}$$

where obviously we have $\sum_{ijklmnl'} \theta_{ijklmnl'} = 1$.

Then we have

$$\begin{aligned} D_X(X) &= \sum_{ijk} \left[\sum_{lmnl'} \tilde{C}_{lmn}(X_B)_{il'} X_{l'l} \tilde{Y}_{jm} \tilde{Z}_{kn} \right. \\ &\quad \left. - \mathcal{A}_{ijk} \ln \left(\sum_{lmnl'} \tilde{C}_{lmn}(X_B)_{il'} X_{l'l} \tilde{Y}_{jm} \tilde{Z}_{kn} \right) \right] + c_1 \\ &\leq \sum_{ijklmnl'} \left[\tilde{C}_{lmn}(X_B)_{il'} X_{l'l} \tilde{Y}_{jm} \tilde{Z}_{kn} \right. \\ &\quad \left. - \mathcal{A}_{ijk} \theta_{ijklmnl'} \ln \frac{\sum_{lmnl'} \tilde{C}_{lmn}(X_B)_{il'} X_{l'l} \tilde{Y}_{jm} \tilde{Z}_{kn}}{\theta_{ijklmnl'}} \right] + c_1 \\ &= \sum_{ijklmnl'} \left[\tilde{C}_{lmn}(X_B)_{il'} X_{l'l} \tilde{Y}_{jm} \tilde{Z}_{kn} \right. \\ &\quad \left. - \gamma_{ijklmnl'} \tilde{B}_{ijk} \ln \left(\tilde{C}_{lmn}(X_B)_{il'} X_{l'l} \tilde{Y}_{jm} \tilde{Z}_{kn} \right) \right] + c_2 \\ &= - \sum_{ijklmnl'} \left[\gamma_{ijklmnl'} \tilde{B}_{ijk} \ln \left(\tilde{C}_{lmn}(X_B)_{il'} X_{l'l} \tilde{Y}_{jm} \tilde{Z}_{kn} \right) \right] + c_3 \\ &\doteq Q(X; \tilde{X}), \end{aligned}$$

where c_1 , c_2 , and c_3 are constants irrelevant to X . Note that in the last step of the above derivation, we used the fact that $\sum_{ijklmnl'} \tilde{C}_{lmn}(X_B)_{il'} X_{l'l} \tilde{Y}_{jm} \tilde{Z}_{kn} = \sum_{ijk} \tilde{C}_{lmn}$ because the columns of X all sum to 1.

It can be easily shown that $Q_X(X; \tilde{X})$ is an *auxiliary function* of $D_X(X)$ in the sense that

$$D_X(X) \leq Q_X(X; \tilde{X}), \text{ and} \quad (16)$$

$$D_X(X) = Q_X(X; X). \quad (17)$$

With such an auxiliary function, we can use the following EM-style argument to show that $X^* = \arg \min_X Q(X; \tilde{X})$ actually reduces $D_X(X)$, namely $D_X(X^*) \leq D_X(\tilde{X})$ in guaranteed:

$$\begin{aligned} D_X(\tilde{X}) &= Q_X(\tilde{X}; \tilde{X}) \text{ (by using Equation (17))} \\ &\geq Q_X(X^*; \tilde{X}) \\ &\geq D_X(X^*) \text{ (by using Equation (16))} \end{aligned}$$

So the problem is reduced to minimizing $Q_X(X; \tilde{X})$ with respect to X , under the constraint that all the columns of X sum to ones. We define the Lagrangian

$$L(X, \vec{\lambda}) = Q(X; \tilde{X}) + \vec{\lambda}^T (X^T \vec{1}_L - \vec{1}_{L'}),$$

and by taking its derivative and setting the result to zero, we have

$$\begin{aligned} \frac{\partial L}{\partial X_{l'l}} &= \frac{\tilde{X}_{l'l}}{X_{l'l}} \sum_{ijkmn} \tilde{B}_{ijk} \tilde{C}_{lmn} \tilde{Y}_{jm} \tilde{Z}_{kn} + \lambda_l = 0 \\ \frac{\partial L}{\partial \lambda_l} &= \sum_{l'} X_{l'l} - 1 = 0 \end{aligned}$$

which gives the update rule for X in Theorem 1. \square

Proof for Corollary 1

A simplex space of dimension $p - 1$ shrunk by ϵ , where $p\epsilon < 1$, is defined as

$$\mathbb{S}_\epsilon^{p-1} = \left\{ x \in \mathbb{R}_+^p : x_k \geq \epsilon \text{ and } \sum_{k=1}^p x_k = 1 \right\}.$$

For a given $\zeta \in \mathbb{R}^p$, and $\max\{\zeta_k\} > 0$, we define a projection

$$\mathcal{P}_\epsilon^p \zeta = \arg \min_{\xi \in \mathbb{S}_\epsilon^{p-1}} - \sum_k \zeta_k \ln \xi_k. \quad (18)$$

Lemma 2 *The loss function is defined as*

$$f(x) = - \sum_i a_i \ln \left(\sum_k w_{ik} x_k \right) - (\alpha - 1) \sum_k \ln(x_k),$$

where $a_i \geq 0$, $\sum_i a_i > p(1 - \alpha)$, and $w_{ik} > 0$. For $x \in \mathbb{S}_\epsilon^{p-1}$, where $\epsilon < 1/p$, if

$$\begin{aligned} b_i &= \frac{a_i}{\sum_k w_{ik} x_k}, \\ \zeta_k &= x_k \sum_i b_i w_{ik} + \alpha - 1, \\ \xi &= \mathcal{P}_\epsilon^p \zeta. \end{aligned}$$

then $\xi \in \mathbb{S}_\epsilon^{p-1}$ and $f(\xi) \leq f(x)$.

Proof We introduce an auxiliary function,

$$\begin{aligned} g(z; x) &= - \sum_{ik} b_i w_{ik} x_k \ln(z_k) + \sum_{ik} b_i w_{ik} x_k \ln \left(\frac{x_k}{\sum_k w_{ik} x_k} \right) - (\alpha - 1) \sum_k \ln(z_k) \\ &= - \sum_k \zeta_k \ln(z_k) + \sum_{ik} b_i w_{ik} x_k \ln \left(\frac{x_k}{\sum_k w_{ik} x_k} \right) \end{aligned}$$

We have $g(x; x) = f(x)$ and $g(z; x) \geq f(z)$ for any z , because of convexity of $-\ln(x)$ in the first term.

We have $\max\{\zeta_k\} > 0$, because $\sum_k \zeta_k = \sum_i a_i + p\alpha - p > 0$. We have $\lambda > 0$, because $\lambda = \zeta_k/\xi_k + \gamma_k \geq \zeta_k/\xi_k \geq \max\{\zeta_k\}/\epsilon > 0$. Thus, $\xi = \mathcal{P}_\epsilon^p \zeta$ minimizes $g(z; x)$. Therefore, $f(\xi) \leq g(\xi; x) \leq g(x; x) = f(x)$.

Although we do not have explicit equation for \mathcal{P}_ϵ^p , inspired by [12], we can solve $\mathcal{P}_\epsilon^p \zeta$ efficiently.

The Lagrangian for Eq. (18) is

$$\mathcal{L} = - \sum_k \zeta_k \ln \xi_k + \lambda \left(\sum_k \xi_k - 1 \right) + \sum_k \gamma_k (\epsilon - \xi_k),$$

where $\gamma_k \geq 0$. With KKT condition, we have $-\frac{1}{\xi_k} \zeta_k + \lambda - \gamma_k = 0$, $\sum_k \xi_k = 1$, and $\gamma_k = 0$ if $\xi_k > \epsilon$.

We prove that $\xi_k \geq \xi_l$ if $\zeta_k \geq \zeta_l$. Suppose that $\xi_l > \xi_k$. If $\xi_l > \xi_k > \epsilon$, it is contradicted by $\zeta_l = \lambda \xi_l > \lambda \xi_k = \zeta_k$. If $\xi_l > \xi_k = \epsilon$, it is contradicted by $\zeta_l = \lambda \xi_l > \lambda \epsilon \geq (\lambda - \gamma_k) \xi_k = \zeta_k$.

We look for ω , such that $\xi_k = \epsilon$ iff $\zeta_k \leq \omega$. Thus $\xi_k = \zeta_k/\lambda > \epsilon$ if $\zeta_k > \omega$. Let $\mathbb{A}_\omega = \{k : \zeta_k \leq \omega\}$, then

$$\frac{1}{\lambda} \sum_{k \notin \mathbb{A}_\omega} \zeta_k + |\mathbb{A}_\omega| \epsilon = 1.$$

So $\lambda = \frac{\sum_{k \notin \mathbb{A}_\omega} \zeta_k}{1 - |\mathbb{A}_\omega| \epsilon}$. Because $\xi_\omega = \epsilon$, $\lambda \geq \omega/\xi_\omega$, thus

$$\epsilon \sum_{\zeta_k > \omega} \zeta_k \geq \omega(1 - |\{\zeta_k \leq \omega\}| \epsilon). \quad (19)$$

Since $\max\{\zeta_k\} > 0$ and $\epsilon < 1/p$, we can find a valid ω such that $\omega \neq \max\{\zeta_k\}$. We select the largest ω that satisfies Eq. (19).

With the analysis above, we solve \mathcal{P}_ϵ^p by Algorithm 7.1.

Algorithm 7.1. $\xi = \mathcal{P}_\epsilon^p \zeta$

- 1: $\{\zeta'\}$ is a descendingly sorted list of $\{\zeta\}$.
- 2: $s_0 = 0$
- 3: **for** $\kappa = 1$ to p **do**
- 4: $s_\kappa = s_{\kappa-1} + \zeta'_\kappa$
- 5: **if** $\epsilon s_\kappa \geq (1 - (p - \kappa)\epsilon)\zeta'_\kappa$ **then**
- 6: $\omega = \zeta'_\kappa$
- 7: $\lambda = s_{\kappa-1}/(1 - (p - \kappa + 1)\epsilon)$
- 8: **break.**
- 9: **end if**
- 10: **end for**
- 11: $\xi_k = \begin{cases} \epsilon, & \text{if } \zeta_k \leq \omega, \\ \zeta_k/\lambda, & \text{if } \zeta_k > \omega. \end{cases}$

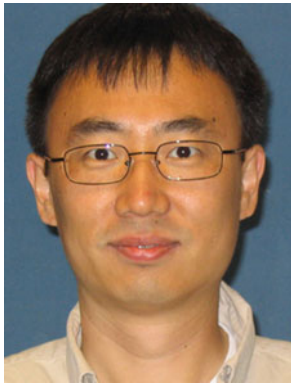
Proof of Corollary 1 Each of \mathcal{C} , \mathcal{X} , \mathcal{Y} , and \mathcal{Z} can be updated using Lemma 2 after certain reshape. Thus we can sequentially minimize the loss function of Eq. (5). \square

References

1. Aussenac-Gilles N, Mothe J (2004) Ontologies as background knowledge to explore document collections. In: RIAO, pp 129–142
2. Banerjee A, Basu S, Merugu S (2007) Multi-way clustering on relation graphs. In: SIAM international conference on data mining
3. Carroll JD, Pruzansky S, Kruskal JB (1980) CANDELINC: a general approach to multi-dimensional analysis of many-way arrays with linear constraints on parameters. *Psychometrika*, 45
4. Chew PA, Bader BW, Kolda TG, Abdelali A (2007) Cross-language information retrieval using PARAFAC2. In: Proceedings of the 13th SIGKDD conference
5. Chi Y, Zhu S (2010) FacetCube: a framework of incorporating prior knowledge into non-negative tensor factorization. In: Proceedings of the 19th CIKM conference
6. Chi Y, Zhu S, Song X, Tatemura J, Tseng BL (2007) Structural and temporal analysis of the blogosphere through community factorization. In: Proceedings of the 13th SIGKDD conference
7. Chi Y, Zhu S, Gong Y, Zhang Y (2008) Probabilistic polyadic factorization and its application to personalized recommendation. In: Proceedings of the 17th CIKM conference
8. Chi Y, Zhu S, Hino K, Gong Y, Zhang Y (2009) iOLAP: a framework for analyzing the internet, social networks, and other networked data. *IEEE Trans Multimedia* 11(3):372–382
9. De Lathauwer L, De Moor B, Vandewalle J (2000) A multilinear singular value decomposition. *SIAM J Matrix Anal Appl* 21(4):1253–1278. doi:[10.1137/S0895479896305696](https://doi.org/10.1137/S0895479896305696)
10. Dhillon IS, Mallela S, Modha DS (2003) Information-theoretic co-clustering. In: Proceedings of the 9th ACM SIGKDD conference
11. Ding C, He X, Simon HD (2005) On the equivalence of nonnegative matrix factorization and spectral clustering. In: SIAM SDM
12. Duchi J, Shalev-Shwartz S, Singer Y, Chandra T (2008) Efficient projections onto the ℓ_1 -ball for learning in high dimensions. In: Proceedings of the 25th international conference on machine learning
13. Fagin R, Lotem A, Naor M (2001) Optimal aggregation algorithms for middleware. In: Proceedings of the twentieth ACM SIGMOD-SIGACT-SIGART symposium on principles of database systems
14. FitzGerald D, Cranitch M, Coyle E (2005) Non-negative tensor factorisation for sound source separation. In: Proceedings of the Irish signals and systems conference
15. Gaussier E, Goutte C (2005) Relation between pls and nmf and implications. In: Proceedings of the 28th SIGIR conference
16. Harshman RA (1970) Foundations of the parafac procedure: models and conditions for an “explanatory” multi-modal factor analysis. UCLA working papers in phonetics, 16
17. Hazan T, Polak S, Shashua A (2005) Sparse image coding using a 3d non-negative tensor factorization. In: Proceedings of the 10th ICCV conference
18. Hofmann T (2001) Unsupervised learning by probabilistic latent semantic analysis. *Mach Learn* 42(1–2):177–196. doi:[10.1023/A:1007617005950](https://doi.org/10.1023/A:1007617005950)
19. Järvelin K, Kekäläinen J (2000) IR evaluation methods for retrieving highly relevant documents. In: Proceedings of the 23rd SIGIR conference
20. Kolda TG, Bader BW (2009) Tensor decompositions and applications. *SIAM Rev* 51(3):455–500
21. Lafferty JD, Zhai C (2001) Document language models, query models, and risk minimization for information retrieval. In: *SIGIR*, pp 111–119
22. Lee DD, Seung HS (2000) Algorithms for non-negative matrix factorization. In: *NIPS*
23. Lin Y-R, Chi Y, Zhu S, Sundaram H, Tseng BL (2008) FacetNet: a framework for analyzing communities and their evolutions in dynamic networks. In: Proceedings of the 17th WWW conference
24. Lin Y-R, Chi Y, Zhu S, Sundaram H, Tseng BL (2009a) Analyzing communities and their evolutions in dynamic social networks. *ACM Trans Knowl Discov Data* 3(2):8:1–8:31. doi:[10.1145/1514888.1514891](https://doi.org/10.1145/1514888.1514891)
25. Lin Y-R, Sun J, Castro P, Konuru R, Sundaram H, Kelliher A (2009b) MetaFac: community discovery via relational hypergraph factorization. In: Proceedings of the 15th ACM SIGKDD international conference on knowledge discovery and data mining
26. Long B, Zhang Z, Yu PS (2007) A probabilistic framework for relational clustering. In: Proceedings of the 13th SIGKDD conference
27. Mørup M, Hansen LK, Arnfred SM (2008) Algorithms for sparse nonnegative Tucker decompositions. *Neural Comput* 20(8):2112–2131
28. Peng W (2009) Equivalence between nonnegative tensor factorization and tensorial probabilistic latent semantic analysis. In: Proceedings of the 32nd international ACM SIGIR conference on research and development in information retrieval
29. Porteous I, Bart E, Welling M (2008) Multi-hdp: a non parametric bayesian model for tensor factorization. In: Proceedings of the 23rd national conference on artificial intelligence

30. Shashua A, Hazan T (2005) Non-negative tensor factorization with applications to statistics and computer vision. In: Proceedings of the 22nd ICML conference
31. Sun J, Faloutsos C, Papadimitriou S, Yu PS (2007) GraphScope: parameter-free mining of large time-evolving graphs. In: Proceedings of the 13th SIGKDD conference
32. Sun J-T, Zeng H-J, Liu H, Lu Y, Chen Z (2005) CubeSVD: a novel approach to personalized web search. In: Proceedings of the 14th WWW conference
33. Tucker LR (1966) Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31
34. Wang F, Li P, Knig A, Wan M (2011) Improving clustering by learning a bi-stochastic data similarity matrix. *Knowl Inf Syst*, pp 1–32
35. Xiong L, Chen X, Huang T-K, Schneider J, Carbonell JG (2010) Temporal collaborative filtering with bayesian probabilistic tensor factorization. In: *SDM*
36. Zaragoza H, Hiemstra D, Tipping ME (2003) Bayesian extension to the language model for ad hoc information retrieval. In: *SIGIR*, pp 4–9
37. Zhang Z-Y, Li T, Ding C (2011) Non-negative tri-factor tensor decomposition with applications. *Knowl Inf Syst*, pp 1–23
38. Zhou D, Zhu S, Yu K, Song X, Tseng BL, Zha H, Giles CL (2008) Learning multiple graphs for document recommendations. In: *WWW*, pp 141–150
39. Zhu S, Yu K, Chi Y, Gong Y (2007) Combining content and link for classification using matrix factorization. In: *SIGIR*

Author Biographies



Yun Chi is currently a Research Staff Member in NEC Laboratories America. He received his Ph.D. degree in Computer Science from University of California, Los Angeles, in 2005. His primary research interests include databases, cloud computing, data mining, machine learning, and information retrieval.



Shenghuo Zhu is currently a Senior Research Staff Member in NEC Laboratories America. He received his Ph.D degree in Computer Science from University of Rochester in 2003. His primary research interests include machine learning, computer vision, data mining and information retrieval.