

Efficient Tensor Decomposition with Application to Topic Modelling

Fengyu Cai

Department of Computer Science
Swiss Federal Institute of Technology Lausanne
Lausanne, Switzerland
fengyu.cai@epfl.ch

Armin Eftekhari

Department of Electrical Engineering
Swiss Federal Institute of Technology Lausanne
Lausanne, Switzerland
armin.eftekhari@gmail.com

Abstract—Nowadays, extracting part of text as information summarization for a collection of documents has become one of popular tasks in the fields of text mining. In this paper, we are going to present a new model for the task of topic extraction. We build up a three-dimensional tensor to describe the collection of documents. And through the mathematical illustration, we decompose the tensor in order to get the mapping relationship between topics and documents with different solver. Based on the decomposed matrix, we rank topics for each document to find its tag.

Index Terms—Topic Selection, Tensor Decomposition

Code Repository:

<https://github.com/TRUMANCFY/SemesterProject2019>

I. INTRODUCTION

Document summarization has become an important research area as a intersection of Information Retrieval (IR), Text Mining (TM), and Natural Language Processing (NLP). It helps a lot in the process of textual repository accessing and data indexing. Based on the huge size of data collection, we put document summerization under the scenario of multiple documents, and the summary of documents will be generated as a form of clusters. Therefore, it is one of key issues about how to present each document and the whole document set. On one hand, the modelling should describe the key features about the documents. On the other hand, the modelling should be general and doable for the further operation on the relationship mapping.

In this paper, we will build a tensor model to describe the document set with the dimentions of documents, terms and potential topics, and decompose it to get the relationship between the document and its topic. Referring to the previously published methods, we generated the topics by clustering the similar sentence. Instead of the traditional high-order Singular Value Decomposition (HOSVD) which has been used in the previous work, we try to solve this using other different solvers. **And in the future, we may check whether the process of clustering can be combined into the solvers. (This direction have not thought about.)**

This paper will be organized as following: in Section II, we will introduce some related work in the field of textual

summarization; in Section III, we will describe the establishment of tensor modelling on the document set; in Section IV, we will discuss how we use different solvers with or without constrains to implement tensor decomposition; Section V will describe the experimental setup and analysis of the result; the final section will contain our conclusion and future work.

II. RELATED WORKED

In the past, the most common method for the document labelling is the computation of the *TF-IDF* score. Even nowadays, around 83% percent of digitized recommender systems are still based on this methodology^{add ref}. It is very straightforward method and shows a stable performance, therefore, we would like to continue to use this as the scoring method. In 2003, Latent Dirichlet allocation (LDA) has been proposed as a generative statistical method for topic labelling. With Dirichlet distribution as prior, it is able to label the document with the probability that it belongs to a specific class. However, the main disadvantage of this method is that the number of classes should be a given parameter. Also, the textual summarization problem has also been transferred to the maximum coverage problem^{add ref}. This methods can extract the sentence to cover as much as key information as possible. As the maximum converge problem is NP-hard, some attempt to find near-optimum solution through a greedy approach. Meanwhile, some other try to apply Integer Linear Programming to find a more precise result.

Tensor is able to well represent high dimensional relationship compactly. And tensor decomposition has been used widely in the bi-relationship extraction from multiple dimensions in many applicable areas. In ^{add ref}, the researchers applied tensor decomposition to summary the comments-oriented documents. They aimed to extract the sentences containing the highly-biased keywords from the collection of comments, based on the rank result from tensor decomposition.

In addition, ^{add ref} provides us the heuristics inspiration about the setup of the tensor. Instead of finding the topics of the document, they relaxed the constrain to find some of the sentences as the summary in a collection of multi-topic documents. They clustered the sentences first and generated

the topics. And then, they built up a tensor with dimensions of terms, topics and documents, and got the ranking of the topics from the decomposition. After that, they selected the sentences out of the clusters based on the order of the topics until reach the maximum length of terms.

In both of two papers mentioned above, they both decomposed the tensor with the method of High Order Singular Value Decomposition (HOSVD). In our project, we intended to use different solver which can provide us more freedom in constrain setting, like Augmented Lagrangian Methods (ALM) or Alternating Direction Method of Multipliers (ADMM). Also, besides the Lagrangian family of methods, we also use get idea from coordinate descent of the non-negative matrix factorization. The previous work [add ref] provides us the existing mathematical tools to simplify our work.

III. OUR METHODS

A. Problem Setting

The main goal of our task to label the document with the topic. Without knowing the topic set, we generate topics first from the set of sentences, based on the assumption that the high dimensional topics generated from sentences is likely to be similar with the topics generated from the documents. After that, we match the document to its closest generated topic, which is most semantically related.

B. Overviews

1) *Pre-processing*: Our pre-processing method follows the standard procedure of the textual data pre-processing, including tokenization, stop-word removal, and stemming.

2) *Generate Topics*: We split the documents into sentences. And we clustered the sentences to find out the potential topics and take them as the document topics.

3) *Tensor Setup*: Given the generated topics, we built up a three-dimensional tensor representing terms, topics and documents.

4) *Tensor Decomposition*: After the setup of the tensor, we get the relationship between topics and terms using different solvers, including the coordinate descent methods and the constrained optimization methods like augmented Lagrangian Methods.

C. Topic Generation

Just Current Methods, will be improved later

Generally speaking, we have a basic assumption that the topics of document is very similar to the topics generated from the sentence. Therefore, we firstly tokenized the document, removed the stopwords, and stemmed the words to the original form. Based on the pre-set number of topics, we use the Latent Dirichlet allocation to generate the topics from the collection of reformed sentences. And also, in this process, we can also get the collections of sentences belonging to the specific topic(s).

D. Tensor Representation Model

Our tensor contains three dimensions. The first dimension of tensor represents the set of terms $T_1...T_m$, the second dimension represents the cluster of topics $C_1...C_k$, and the third dimension presents the collection of documents $D_1...D_n$. Therefore, each entry $\mathbf{T}_{i,j,k}$ in the tensor contains the information of i -th term related to j -th topic in the k -th document. Based on the last step *Topic Generation*, we can find out that some words belong to the specific topics, therefore, it is an advantage that we only need to consider these kind of entries instead of the whole tensor.

$$\mathbf{T}_{i,j,k} = \begin{cases} tf - idf_{i,j,k} & T_i \in C_j \text{ and } T_i \in D_k \\ 0 & \text{otherwise} \end{cases}$$

There will be a figure to show the tensor representative

E. Tensor Decomposition Methods

After getting the tensor, we tried to decompose it to get the relational mapping between documents and topics. Generally speaking, we categorized the methods to unconstrained and constrained problem.

The process of tensor decomposition could be interpreted to find three matrices U, V and W, such that

$$\mathbf{T}_{i,j,k} \approx U_{i,r} \times V_{j,r} \times W_{k,r} \quad (1)$$

Therefore, we can let

$$\mathcal{L}(\mathbf{T}; \mathbf{U}, \mathbf{V}, \mathbf{W}) = \frac{1}{2} \sum_{i,j,k}^{I,J,K} (\mathbf{T}_{ijk} - \sum_{r=1}^R U_{ir} \times V_{jr} \times W_{kr}) \quad (2)$$

to be the loss function. Here, R denote the rank of the reconstructed tensor, I denotes the total number of documents, J denotes the total number of topics, and K denotes the total number of terms. Therefore, U contains the information about documents, V contains the information about topics, and W will be related to terms. In order to get the matrix representing the relationship between documents and topics, we set the rank R equal to the number of topics. Our task will be transfer to

$$\min_{\mathbf{U}, \mathbf{V}, \mathbf{W}} \mathcal{L}(\mathbf{T}; \mathbf{U}, \mathbf{V}, \mathbf{W}) \quad (3)$$

and we will mainly focus on our target U and add constrain on it in the further.

1) Alternative Gradient Descent:

As we know, the gradient descent method is described as followed:

$$u^{t+1} = u^t - \alpha \frac{\partial \mathcal{L}}{\partial u}(u^t; x^t) \quad (4)$$

if the loss function is $\mathcal{L}(x; u)$, where t is the step of iteration and α is the learning rate.

The main difficulty of this problem is to get the gradient. The paper (add ref) provides the mathematical tools to derive the gradient. Firstly, we define a matrix $\mathbf{T}(\cdot, \mathbf{V}, \mathbf{W}) \in \text{Real}^{I \times R}$, from \mathbf{T} , \mathbf{V} and \mathbf{W} :

$$\mathbf{T}(\cdot, \mathbf{V}, \mathbf{W})_{ir} := \sum_{jk} \mathbf{T}_{ijk} V_{jr} W_{kr} \quad (5)$$

Therefore, according to **Theorem 1** mentioned in the paper, CP gradient of loss function can be written as:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{U}}(\mathbf{T}; \mathbf{U}, \mathbf{V}, \mathbf{W}) = -\mathbf{T}(\cdot, \mathbf{V}, \mathbf{W}) + \mathbf{U}\Gamma((\mathbf{V}), \mathbf{W}) \quad (6)$$

where $\Gamma(A, B)$ is the Hadamard product of $A^T A$ and $B^T B$, i.e., $\Gamma(A, B)_{ij} = (A^T A)_{ij} (B^T B)_{ij}$. And the formulas of $\partial \mathcal{L} / \partial \mathbf{V}$ and $\partial \mathcal{L} / \partial \mathbf{W}$.

Therefore, we can update \mathbf{U} , \mathbf{V} and \mathbf{W} iteratively, setting two of them as constant variable and update the other. **also cyclic variant version: use \mathbf{U}^{t+1} to update \mathbf{V}^{t+1} and \mathbf{W}^{t+1}**

$$\mathbf{U}^{t+1} = \mathbf{U}^t - \alpha \frac{\partial \mathcal{L}}{\partial \mathbf{U}}(\mathbf{T}; \mathbf{U}^t, \mathbf{V}^t, \mathbf{W}^t) \quad (7)$$

$$\mathbf{V}^{t+1} = \mathbf{V}^t - \alpha \frac{\partial \mathcal{L}}{\partial \mathbf{V}}(\mathbf{T}; \mathbf{U}^t, \mathbf{V}^t, \mathbf{W}^t) \quad (8)$$

$$\mathbf{W}^{t+1} = \mathbf{W}^t - \alpha \frac{\partial \mathcal{L}}{\partial \mathbf{W}}(\mathbf{T}; \mathbf{U}^t, \mathbf{V}^t, \mathbf{W}^t) \quad (9)$$

2) Alternative Gradient Descent with regularization:

However, as the optimization is non-convex and multi-linear, it may be stuck in ill-condition. There exit a lot of local minimum due to the scale of three values. For example, $\mathcal{L}(\mathbf{T}; \mathbf{U}, \mathbf{V}, \mathbf{W}) = \mathcal{L}(\mathbf{T}; a\mathbf{U}, b\mathbf{V}, \mathbf{W}/ab)$

Therefore, we add the regularization terms to the loss function:

$$\mathcal{L}_\rho(\mathbf{T}; \mathbf{U}, \mathbf{V}, \mathbf{W}) = \mathcal{L}(\mathbf{T}; \mathbf{U}, \mathbf{V}, \mathbf{W}) + \frac{\rho}{2} (\|\mathbf{U}\|^2 + \|\mathbf{V}\|^2 + \|\mathbf{W}\|^2) \quad (10)$$

The gradient of regularized loss function is very similar to the original one:

$$\frac{\partial \mathcal{L}_\rho}{\partial \mathbf{U}}(\mathbf{T}; \mathbf{U}, \mathbf{V}, \mathbf{W}) = -\mathbf{T}(\cdot, \mathbf{V}, \mathbf{W}) + \mathbf{U}\Gamma_\rho((\mathbf{V}), \mathbf{W}) \quad (11)$$

where $\Gamma_\rho(A, B)_{ij} = \Gamma(A, B)_{ij} + \rho I_{ij}$. Also $\partial \mathcal{L} / \partial \mathbf{W}$ and $\partial \mathcal{L} / \partial \mathbf{V}$ are in the similar formulas. And the process of iterative gradient descent can be referred to last non-regularized part.

3) Second-order Gradient Descent:

Compared with the first-order gradient descent, the second-order is able to converge in fewer iterations. **may need ref** The format is defined as:

$$\mathbf{u}^{t+1} = \mathbf{u}^t - \alpha \mathcal{H}^t \frac{\partial \mathcal{L}}{\partial \mathbf{u}}(\mathbf{u}^t; \mathbf{x}^t) \quad (12)$$

where \mathcal{H} is the approximate inverse Hessian matrix.

Based on the derivative **(may need more explanation and references, when increasing the rank, it does not work)**, we could update \mathbf{U} in the following formula:

$$\mathbf{U}^{t+1} = \mathbf{U}^t - \alpha \frac{\partial \mathcal{L}}{\partial \mathbf{U}} \Gamma_\rho(\mathbf{V}^t, \mathbf{W}^t)^{-1} \quad (13)$$

$$= \mathbf{U}^t (1 - \alpha) + \alpha \mathbf{T}(\cdot, \mathbf{V}^t, \mathbf{W}^t) \Gamma_\rho(\mathbf{V}^t, \mathbf{W}^t)^{-1} \quad (14)$$

4) Alternative Directional Methods of Multiplier:

Till now, three methods mentioned above are all non-constrained. Therefore, we can add some constrain to our optimization problem. For example, we could let the entry \mathbf{U}_{ir} to be the probability that the topic of i -th document is the r -th topic. Therefore, the following constrain should be satisfied for i -th row in \mathbf{U} :

$$\sum_{r=1}^R \mathbf{U}_{ir} = 1 \quad (15)$$

If we vectorize the equation above, we set up the constrain:

$$\mathbf{U}\mathbf{1} = \mathbf{1} \quad (16)$$

where $\mathbf{1}$ is a vector filled by one with suitable dimension.

And then, we can introduce Lagrangian multiplier \mathcal{Y} to combine the original loss function and the constrain. Also, in order to guarantee the smoothness of the dual function, we can add the penalty term $\frac{\mu}{2} \|\mathbf{U}\mathbf{1} - \mathbf{1}\|^2$ into the new loss function. Therefore, considering two additional terms, we define the loss function of ALM as followed:

$$\mathcal{L}_\mu(\mathbf{T}; \mathbf{U}, \mathbf{V}, \mathbf{W}, \mathcal{Y}) = \mathcal{L}(\mathbf{T}; \mathbf{U}, \mathbf{V}, \mathbf{W}) + \langle \mathcal{Y}, \mathbf{U}\mathbf{1} - \mathbf{1} \rangle + \frac{\mu}{2} \|\mathbf{U}\mathbf{1} - \mathbf{1}\|^2 \quad (17)$$

According to the procedure of Lagragian method, given a \mathcal{Y} , we need to find out the optimized \mathbf{U} :

$$\mathbf{U}^{t+1} = \min_{\mathbf{U}} \mathcal{L}_\mu(\mathbf{T}; \mathbf{U}^t, \mathbf{V}^t, \mathbf{W}^t, \mathcal{Y}^t) \quad (18)$$

However, it is very difficult to get \mathbf{U}^{t+1} directly, therefore, we need to linearize it:

$$\mathbf{U}^{t+1} = \min_{\mathbf{U}} \mathcal{L}_\mu(\mathbf{T}; \mathbf{U}^t, \mathbf{V}^t, \mathbf{W}^t, \mathcal{Y}^t) \quad (19)$$

$$\leq \min_{\mathbf{U}} \mathcal{L}(\mathbf{T}; \mathbf{U}^t, \mathbf{V}^t, \mathbf{W}^t) + \frac{\mu}{2} \|\mathbf{U}^t \mathbf{1} - \mathbf{1}\|^2 \quad (20)$$

$$+ \mu \langle (\mathbf{U}^t \mathbf{1} - \mathbf{1}) \mathbf{1}^T, \mathbf{U} - \mathbf{U}^t \rangle \quad (21)$$

$$+ \langle \mathcal{Y}^t, \mathbf{U} - \mathbf{U}^t \rangle + \frac{\mu}{2} \|\mathbf{U} - \mathbf{U}^t\|^2 \quad (22)$$

$$= \mathbf{F}_{\mu, \mathcal{Y}^t}(\mathbf{U}) \quad (23)$$

In order to find the optimized \mathbf{U} in each step, we let $\frac{\partial \mathbf{F}}{\partial \mathbf{U}} = 0$. In the other words,

$$\frac{\partial \mathbf{F}}{\partial \mathbf{U}} = -\mathbf{T}(\cdot, \mathbf{U}^t, \mathbf{V}^t) + \mathbf{U}\Gamma(\mathbf{V}^t, \mathbf{W}^t) \quad (24)$$

$$+ \mu(\mathbf{U} - \mathbf{U}^t) + \mu(\mathbf{U}^t \mathbf{1} - \mathbf{1})\mathbf{1}^T + \mathcal{Y} = 0 \quad (25)$$

Therefore, we can update \mathbf{U} as followed:

$$\mathbf{U}^{t+1} = \frac{1}{\mu}(\mathbf{T}(\cdot, \mathbf{U}^t, \mathbf{V}^t) - \mathcal{Y} - \mu(\mathbf{U}^t \mathbf{1} - \mathbf{1})\mathbf{1}^T + \mu\mathbf{U}^t)(\mathbf{I} + \Gamma(\mathbf{V}^t, \mathbf{W}^t))^{-1} \quad (26)$$

And then, like general Lagrangian methods, we do the gradient ascent on \mathcal{Y} to update it.

5) *Proximal Operator*: In order to keep the sparsity of the document-topic matrix, we can add the nuclear-norm to the loss function as constrain:

$$\mathcal{L}_{\lambda\|\cdot\|_*}(\mathbf{T}; \mathbf{U}, \mathbf{V}, \mathbf{W}) = \mathcal{L}(\mathbf{T}; \mathbf{U}, \mathbf{V}, \mathbf{W}) + \frac{\lambda}{2}\|\mathbf{U}\|_* \quad (27)$$

Also, it can be added with other constrained term. However, the main bottleneck is that the loss function will not be smooth anymore. Subgradient can be applied, but the convergence rate will be largely affected and reduced to $\mathcal{O}(\frac{1}{k})$. [add ref about the convergence rate](#) Therefore, the proximal operate of nuclear norm can solve this issue. The operator will find out a close-form solution which is close to the input, but try to satisfy the constrain.

$$\mathbf{Prox}_{\lambda\|\cdot\|_*}(\mathbf{A}) = \arg \min_{\mathbf{X}} \frac{1}{2}\|\mathbf{A} - \mathbf{X}\|_F^2 + \frac{\lambda}{2}\|\mathbf{X}\|_* \quad (28)$$

After a series of derivation [maybe add in the future](#), we could take advantage of proximal operator of L1-norm, which is easy to derive:

$$\mathbf{Prox}_{\lambda\|\cdot\|_1}(\mathbf{a}) = \arg \min_{\mathbf{x}} \frac{1}{2}\|\mathbf{a} - \mathbf{x}\|_2^2 + \frac{\lambda}{2}\|\mathbf{x}\|_1 \quad (29)$$

$$= \text{sign}(x) \max(x - \lambda, 0) \quad (30)$$

Therefore, we can do Singular Value Decomposition(SVD) on \mathbf{A} , i.e. $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T$, where \mathbf{U} and \mathbf{V} are the orthogonal basis, and \mathbf{S} is the diagonal matrix containing the singular values. As the nuclear norm of the diagonal matrix with singular values is equivalent to the L1-norm of the vector containing the singular values, the