



UNIVERSIDAD NACIONAL DE COLOMBIA
SEDE BOGOTÁ

FACULTAD DE INGENIERÍA
Vicedecanatura Académica
POSGRADOS

PROPOSAL SUBMISSION

DOCTORAL THESIS: ☒ MASTER THESIS: ☐
MASTER FINAL WORK: ☐ SPECIALIZATION FINAL WORK: ☐

1. **BIDDER:** Robinson Andrés Jaque Pirabán **ID:** 80190790
2. **PROGRAM:** Philosophy Doctoral in Computer Science and Systems Engineering
3. **ADVISOR:** Fabio Augusto González Osorio
DEPARTMENT: Computer Science and Industrial Engineering
4. **TITLE:** Kernel Tensor Factorization
5. **AREA:** Computer Science
6. **LINE OF RESEARCH:** Machine Learning
7. COMMENTARY WITH ADVISOR APROVAL

8. BIDDER SIGNATURE

9. SIGNATURE OF ADVISOR

Contents

1	Introduction	3
1.1	Basics of tensors	3
1.1.1	Unfolding and Folding Tensors	4
2	Tensor Decomposition	4
2.1	Tensor Factorization Methods	6
2.1.1	Canonical Polyadic Decomposition / PARAFAC	6
2.1.2	TUCKER Decomposition	7
2.1.3	Non-negative Tensor Factorization	8
2.1.4	Other methods	9
2.2	Kernel methods	9
3	Kernel Non-negative Matrix Factorization	10
3.1	General problems	10
3.2	Tensor completion	10
4	Problem Statement	11
5	Goals	12
6	Justification	13

1 Introduction

1.1 Basics of tensors

Tensors are multidimensional arrays, i.e. an N -way or N -order tensor is an element of tensor product of N vector spaces, each of which has its own coordinate system. A first order tensor is a vector, a second order tensor is a matrix, tensors of higher order are called high-order tensors. The order (ways or modes) of a tensor is the number of dimensions. Figure 1 represents a 3-order tensor $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$.

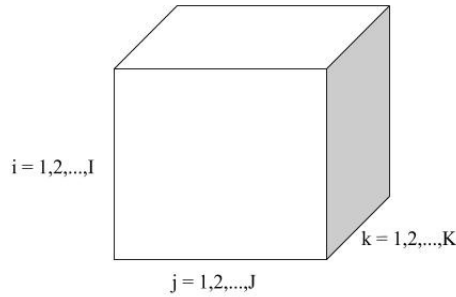


Figure 1: Third-order tensor

Fibers are defined by fixing every index by one. In a third-order tensor a column is a mode-1 fiber, denoted by $x_{:jk}$; a row is a mode-2 fiber, denoted by $x_{i:k}$; while a tube is a mode-3 fiber, denoted by $x_{i:k}$. 2 shows a fibers representation in 3rd-order tensor.

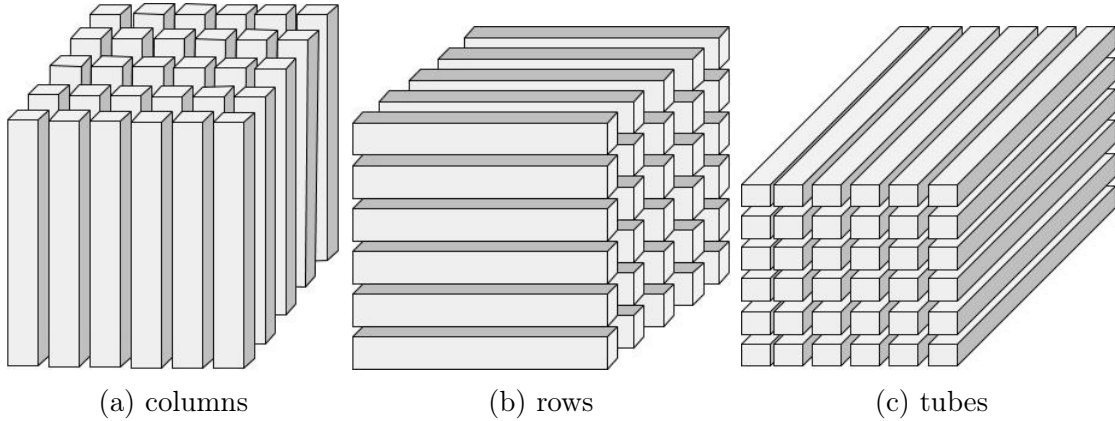


Figure 2: 3rd-order tensor fibers

Slices are two-dimensional sections of a tensor defined by fixing two indexes. For instance, slices of 3rd-order tensor \mathcal{X} are denoted by $X_{i:}$ (horizontal), $X_{:j}$ (lateral) and $X_{::k}$ (frontal) and we illustrate them in figure 3.

The *norm* of a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ is analogous to the matrix Frobenius norm,

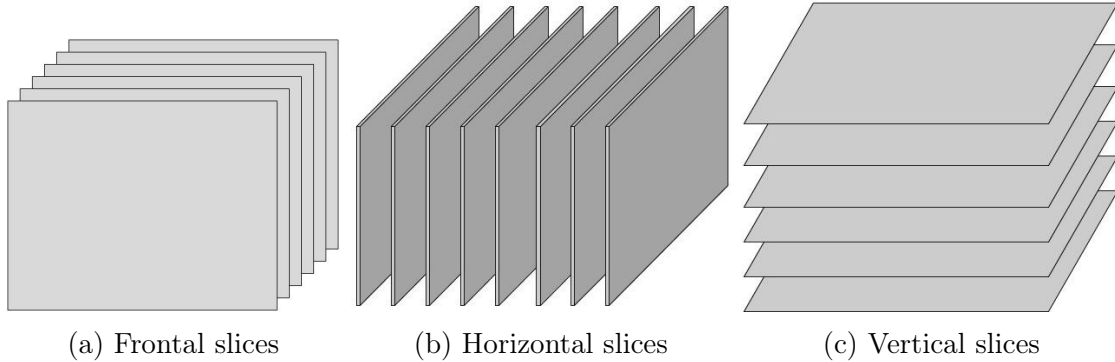


Figure 3: 3rd-order tensor slices

i.e.

$$\|\mathcal{X}\| = \sqrt{\sum_{i_1=1}^{I_1} \sum_{i_2=1}^{I_2} \cdots \sum_{i_N=1}^{I_N} x_{i_1 i_2 \dots i_N}^2} \quad (1)$$

\mathcal{X} is a *Rank-one* tensor if it is equal to the outer product of N vectors, i.e.,

$$\mathcal{X} = a^{(1)} \otimes a^{(2)} \otimes \dots \otimes a^{(N)}$$

1.1.1 Unfolding and Folding Tensors

Unfolding is the process of *matricization* of a tensor. In other words, elements of a tensors are sorted to assemble a matrix. The mode- k unfolding of a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ is denoted by $X_{(k)} \in \mathbb{R}^{I_1 \times \prod_{k' \neq k} I_{k'}}$ and arranges the mode- k tensor fibers as columns of resulting matrix. In addition, Kolda [?] presents a more general procedures of unfolding

Ding and Wei [?] present a fast algorithm for Hankel tensor-vector products.

[?] A method of fast linear transform algorithm synthesis for an arbitrary tensor

2 Tensor Decomposition

Tensor decomposition originated with Hitchcock in 1927 [?], and the the multi-way model is attributed to Cattell in 1944 [?].

Tensor works had attention in 60s with Tucker ([?], [?], [?]) and Carroll and Chang [?] and Harshman in 1970 [?] with applications in psychometrics. In 1981 Appelhof and Davidson [?] used tensor decomposition in chemometrics which have been an popular field of application of tensor decomposition since then.

In last twenty years tensor decomposition applications have expanded to many fields such as signal processing, numerical linear algebra, computer vision, numerical analysis, neuroscience, data mining, graph analysis. Figure 4 shows seminal papers which opened broad application fields to tensor decomposition, plot also shows the amount of documents published about these works according to Scopus.

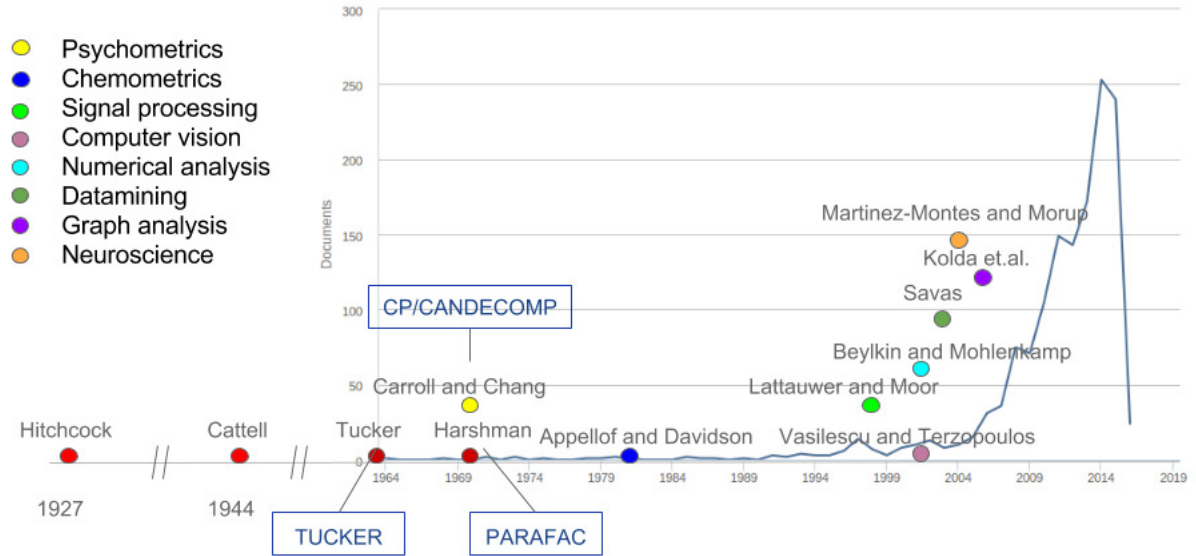


Figure 4: Time-line of tensor decomposition

Kolda [?], Acar [?] and [?] present an exhaustive and detailed review of fundamental decomposition methods and applications. Furthermore, [?] presents tensor properties as extension of structural properties of matrices. On the other hand, Fanaee and Gama [?] introduce an interdisciplinary survey about tensor-based anomaly detection.

In following sections we explain some of basic methods to tensor decomposition which have been inspiration to many others methods proposed. Also, we summarize recently works and approaches of tensor decomposition on different application fields.

Chemometrics [?] A computationally efficient technique for the solution of multi-dimensional PBMs of granulation via tensor decomposition

Bioinformatics In [?] and [?] a multimodal problem is addressed, the authors formulate data fusion as a coupled matrix and tensor factorization problem and discuss its extension to a structure-revealing data fusion model in metabolomics.

Image processing [?] NNTF for facial expression recognition

Machine Learning [?] Tensor decompositions for learning latent variable models

Text mining [?] This paper describes a method for automatic detection of semantic relations between concept nodes of a networked ontological knowledge base by analyzing matrices of semantic-syntactic valences of words. These matrices are obtained by means of nonnegative factorization of tensors of syntactic compatibility of words.

Numerical analysis [?] use of the sum-factorization for the calculation of the integrals arising in Galerkin isogeometric analysis. While introducing very little change

in an isogeometric code based on element-by-element quadrature and assembling, the sum-factorization approach, taking advantage of the tensor-product structure of splines or NURBS shape functions, significantly reduces the quadrature computational cost.

[?] Fast iterative solution of the Bethe-Salpeter eigenvalue problem using low-rank and QTT tensor approximation.

[?] Blind identification of a second order Volterra-Hammerstein series using cumulant cubic tensor analysis.

Neuroscience [?] Decomposition of brain diffusion imaging data uncovers latent schizophrenias with distinct patterns of white matter anisotropy, using NNTF to clustering.

Signal processing Cichocki et.al. [?] sum up tensor decomposition approaches for signal processing problems.

Barker and Virtanen [?] deal with monaural sound source separation problem using NNTF of modulation spectrograms.

[?] Necessity to manually assign the NTF components to audio sources in order to be able to enforce prior information on the sources during the estimation process, Automatic Allocation of NTF Components for User-Guided Audio Source Separation

[?] propose a shifted 2D non-negative tensor factorisation algorithm which extends non-negative matrix factor 2D deconvolution to the multi-channel case. The use of this algorithm for multi-channel sound source separation of pitched instruments is demonstrated.

Other applications [?] Discovering and Characterizing Mobility Patterns in Urban Spaces: A Study of Manhattan Taxi Data. by using non-negative tensor factorization (NTF), we are able to cluster human behavior based on spatio-temporal dimensions. Second, for understanding these clusters, we propose to use HypTrails, a Bayesian approach for expressing and comparing hypotheses about human trails.

[?] NTF factorization for household electrical seasonal consumption disaggregation

2.1 Tensor Factorization Methods

Main

2.1.1 Canonica Polyadic Decomposition / PARAFAC

Canonica Polyadic (CP) decomposition (Kolda,2009) (Mocks, cite 166 Kolda), CAN-DECOMP(Carroll and Chang, cite 70 Kolda) or PARAFAC (Harshman, cite 90 Kolda) decompose a tensor as a finite sum of rank-one tensors. For instance, given a third order tensor $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$, CP decomposition express it as

$$\mathcal{X} \approx \sum_{r=1}^R a_r \otimes b_r \otimes c_r \quad (2)$$

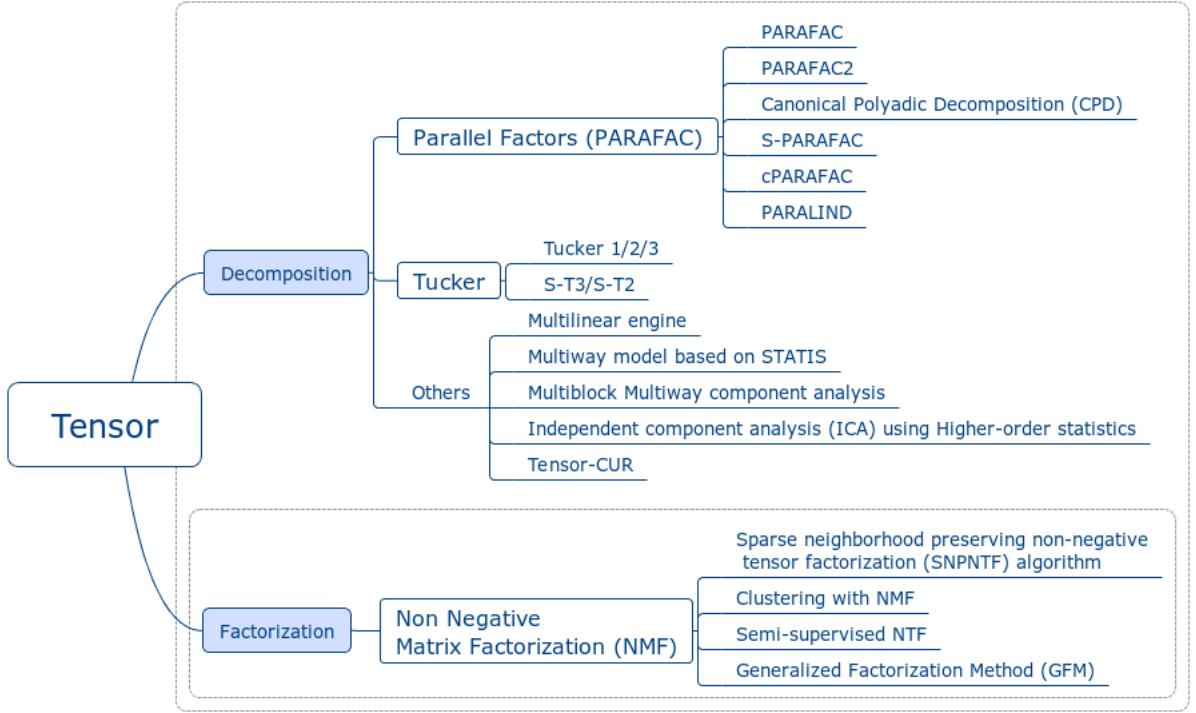


Figure 5: Tensor decomposition methods

where $a_r \in \mathbb{R}^I$, $b_r \in \mathbb{R}^J$, $c_r \in \mathbb{R}^K$ and R is a positive integer.

Domanov [?] shows relaxed uniqueness conditions and algebraic algorithm for Canonical polyadic decomposition, as well as a reduction to generalized eigenvalue decomposition [?] and uniqueness properties [?] of third-order tensors.

2.1.2 TUCKER Decomposition

Tucker decomposition was introduced by Tucker (Tucker, 1963, 1966). It is also named N-mode PCA (kapteyn, 1986, cite 113 Kolda), High-order SVD (HOSVD) (De Lathauwer, 2000, cite 63 Kolda) or N-mode SVD (Vasilescu, 2002, cite 229 Kolda).

The Tucker decomposition is a form of higher-order PCA (Kolda, 2009). It decomposes a tensor into a core tensor \mathcal{G} multiplied by a matrix along each mode. For instance, given a third order tensor $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$, Tucker decomposition express it as

$$\mathcal{X} \approx \mathcal{G} \times_1 A \times_2 B \times_3 C \quad (3)$$

Where \times_k is the mode- k product, $A \in \mathbb{R}^{I \times P}$, $B \in \mathbb{R}^{J \times Q}$, $C \in \mathbb{R}^{K \times R}$ are the factor matrices (usually orthogonals) and can be interpreted as the principal components for each mode. $\mathcal{G} \in \mathbb{R}^{P \times Q \times R}$ is the core tensor and its entries show the interactions between the different components.

2.1.3 Non-negative Tensor Factorization

Non-negative Matrix Factorization The general problem of non-negative matrix factorization (NMF) is to decompose a matrix $X \in \mathbb{R}_{\geq 0}^{n \times l}$ into two matrix factors: basis $W \in \mathbb{R}_{\geq 0}^{n \times k}$ and coefficients $H \in \mathbb{R}_{\geq 0}^{k \times l}$, i.e.

$$X \cong WH \quad (4)$$

The factorization problem can be seen as an optimization problem:

$$\min_{W, H} d(X, WH) \quad (5)$$

where $d(\cdot)$ is a distance or divergence function and the problem could have different types of restrictions. For instance, if $d(\cdot)$ is the Euclidean Distance and there are not restrictions, the problem is solved by finding the SVD; if X , W and H are restricted to be positive, then the problem is solved by NMF. An comprehensive survey of NMF variants and algorithms is found in (Wang and Zhang, 2013)

One NMF approach is Symmetric-NMF [?], (SNMF) which produces a factorization:

$$(X_{l \times n}^T X_{n \times l}) = H_{l \times k} H_{k \times l}^T \quad (6)$$

An important characteristic of this version of NMF is that it is amenable to be used as a kernel method. This is discussed in the next subsection.

Non-negative Tensor Factorization An natural extension of nonnegative matrix factorization with high-order arrays is nonnegative n-dimensional tensor factorization (n-NTF). This kind of generalization is indeed not trivial since NTF possesses many new properties varying from NMF (cites [117], [118] of (Wang and Zhang, 2013))

(cp Wang and Zhang, 2013) First, the data to be processed in NMF are vectors in essence. However, in some applications the original data may not be vectors, and the vectorization might result in some undesirable problems. For instance, the vectorization of image data, which is two dimensional, will lose the local spatial and structural information. Second, one of the core concerns in NMF is the uniqueness issue, and to remedy the ill-posedness some strong constraints have to be imposed. Nevertheless, tensor factorization will be unique under only some weak conditions. Besides, the uniqueness of the solution will be enhanced as the tensor order increases.

There are generally two types of NTF model—NTD [119] and more restricted NPARAFAC [117], whose main difference lies in the core factor tensor. As for the solution, there are some feasible approaches. For example, NTF can be restated as regular NMF by matricizing the array [116], [119]. Or the alternating iteration method can be utilized directly on the outer product definition of tensors [117], [118], [23]. Similarly, SED, GKLD and other forms of divergence can also be used as the objective functions [23], [120], [121]. And some specific update models can adopt the existing conclusions in NMF. For thorough understanding one may refer to [9],

[122]. What must be scrutinized here is that the convergence of these algorithms is not guaranteed by the simple generation from matrix to tensor forms in itself.

What's more, the concepts in Constrained NMF can also be incorporated in NTF, such as sparse NTF [123], [124], [125], discriminant NTF [126], NTF on manifold [127], and the like.

2.1.4 Other methods

The tensor decomposition addressed in [?] may be seen as a generalization of Singular Value Decomposition of matrices. They consider general multilinear and multihomogeneous tensors. Then, they show how to reduce the problem to a truncated moment matrix problem and give a new criterion for flat extension of Quasi-Hankel matrices.

2.2 Kernel methods

In contrast with traditional learning techniques, kernel methods do not need a vectorial representation of data. Instead, they use a kernel function. Therefore, kernel methods are naturally applied to unstructured, or complex structured, data such as texts, strings, trees and images [?].

Informally, a kernel function measures the similarity of two objects. Formally, a kernel function, $k : X \times X \rightarrow \mathbb{R}$, maps pairs (x, y) of objects in a set X , the problem space, to the reals. A kernel function implicitly generates a map, $\Phi : X \rightarrow F$, where F corresponds to a Hilbert space called the feature space. The dot product in F is calculated by k , specifically $k(x, y) = \langle \Phi(x), \Phi(y) \rangle_F$. Given an appropriate kernel function, complex patterns in the problem space may correspond to simpler patterns in the feature space. For instance, non-linear patterns in the problem space may correspond to linear patterns in the feature space.

Both k -means and SNMF have kernelized versions, which receive as input a kernel matrix instead of a set of sample represented by feature vectors. The kernel version of k -means is called, unsurprisingly, kernel k -means (KKM). In the case of SNMF, the kernelized version works as follows.

SNMF starts with an initial estimation of the matrix factor H and iteratively update it using the updating equation:

$$H_{i,k} = H_{i,k} (1 - \beta + \beta \frac{((X^T X)H)_{i,k}}{(HH^T H)_{i,k}})$$

The kernel version of the algorithm is obtained by using a kernel matrix K instead of the expression $(X^T X)$, where K is an $l \times l$ matrix with $K_{i,j} = k(x_i, x_j)$. There are different types of kernels some of them general and some of them specifically defined for different types of data. The most popular general kernels are the linear kernel

$$k(x, y) = \langle x, y \rangle, \tag{7}$$

the polynomial kernel

$$k(x, y) = p(\langle x, y \rangle),$$

where $p(\cdot)$ is a polynomial with positive coefficients, and the Gaussian (or RBF) kernel

$$k(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}}. \quad (8)$$

The cluster centroids estimated by the kernel versions of both algorithms are in the feature space and correspond to the points $C_j = \frac{1}{n} \sum_{x_i \in C_j} \Phi(x_i)$. However, we are interested on the pre-image in the original space of this centroids, i.e., points \hat{C}_j such that $\Phi(\hat{C}_j) = C_j$. However, it is possible that a exact pre-image may not even exist, so we look for the \hat{C}_j that minimizes the following objective function: $\min_{\hat{C}_j} \|\hat{C}_j - C_j\|^2$. According to Kwok et al. [?], the optimum C_j can be found by iterating the following fixed-point formula:

$$\hat{C}_j^{t+1} = \frac{\sum_{i=1}^N \exp\left(\frac{-\|\hat{C}_j^t - x_i\|}{s}\right) x_i}{\sum_{i=1}^N \exp\left(\frac{-\|\hat{C}_j^t - x_i\|}{s}\right)} \quad (9)$$

3 Kernel Non-negative Matrix Factorization

Kernel Non-negative Matrix Factorization (KNMF) can be naturally derivated of convex NMF (insert cites 92, 98 and 120 from Cichocki book). Given a kernel function $\phi : x \in X \rightarrow \phi(x) \in F$, mapping for N elements $\phi(X) = [\phi(x_1), \dots, \phi(x_N)]$. Then, KNMF can be defined as

$$\phi(X) \cong \phi(X)WH^T \quad (10)$$

Therefore, the cost function to minimize is

$$\|\phi(X) - \phi(X)WH^T\|_F^2 = \text{tr}(K) - 2\text{tr}(H^T KW) + \text{tr}(W^T KWH^T H) \quad (11)$$

Where kernel $K = \phi^T(X)\phi(X)$

3.1 General problems

Usually, tensor factorization address the folowing problems independent of the application: blind source separation, tensor completion.

3.2 Tensor completion

In tensor completion, an given incomplete tensor is given, i.e. some of its entries are missing and we should complete that. Following Ji Lu et.al [?] notation, low rank matrix completion is noted as

$$\begin{aligned} & \min_X \text{rank}(X) \\ & \text{s.t. } X_\Omega = M_\Omega \end{aligned} \quad (12)$$

where Ω is an index set, then X_Ω is coping entries of X in the indexes Ω and missed entries $\hat{\Omega}$ would be 0

The missing entries in X are determined in order to minimize the matrix X rank. i.e. a non convex optimization problem since rank is nonconvex.

Frequently, trace norm (or nuclear norm) $\|\cdot\|_*$ is used to approximate the rank of matrices.

Trace norm is the tightest convex envelop for the matrices rank.

$$\begin{aligned} \min_X & \|X\|_* \\ \text{s.t. } & X_\Omega = M_\Omega \end{aligned} \quad (13)$$

Since tensor is a generalization of the matrix concept, we generalize the optimization problem as

$$\begin{aligned} \min_{\mathcal{X}} & \|\mathcal{X}\|_* \\ \text{s.t. } & \mathcal{X}_\Omega = \mathcal{T}_\Omega \end{aligned} \quad (14)$$

Where \mathcal{X} and \mathcal{T} are n -order tensors with identical size.

Acar et.al [?] presents a scalable tensor factorization method to deal with completion problem using PARAFAC method. Cao [?] propose a new tensor completion model via folded-concave penalty for estimating missing values in tensor data. To solve the resulting nonconvex optimization problem, we develop a local linear approximation augmented Lagrange multiplier (LLA-ALM) algorithm which combines a two-step LLA strategy to search a local optimum of the proposed model efficiently. They show numerical results in image and video data sets and compare with nuclear norm penalization method in order of demonstrate its advantage in terms of the accuracy and robustness.

[?] propose a method to deal with the completion problem when the number of missing entries increases, since factorization schemes may overfit the model because of incorrectly predefined ranks, while completion schemes may fail to interpret the model factors. Hence, they present an approach to complete the missing entries and simultaneously capture the underlying model structure that combines a rank minimization technique with Tucker model decomposition. Moreover, as the model structure is implicitly included in the Tucker model, they use factor priors, which are usually known a priori in real-world tensor objects, to characterize the underlying joint-manifold drawn from the model factors.

4 Problem Statement

Matrix and analogous tensor factorization is central to different important tasks in machine learning and information retrieval such as: clustering, latent topic analysis, recommendation, blind source separation, completion, among others. The widespread use of multisensor technology and the emergence of big data sets have highlighted the

limitations of standard flat-view matrix models and the necessity to move toward more versatile data analysis tools (Cichocki, 2015). Conventional methods preprocess multiway data by arranging them into a matrix, which might lose the original multiway structure of the data (Wang and Zhang, 2013). Hence, tensors address multimodal or multiview data, and tensor factorization brings tools to perform common machine learning tasks.

On the other hand, kernel methods are ubiquitous in machine learning, performance of these methods have been broadly demonstrated, and there are evidence between some types of kernels and robustness. However, there are few exploration of kernel tensor factorization approaches.

A satisfactory solution of this general challenge requires to answer some particular **research question**:

- How incorporate kernel methods in tensor factorization in order to deal with multiway data?

An answer to the question derivate in a method which factorize a given tensor in the feature space induced by a Kernel function. Naturally, to inquire about the effects of decompose tensors in an space induced by a kernel function open space to evaluate the proposal method performance as well as its capabilities dealing with multimodal data, scalability and robustness to noise and outliers.

5 Goals

The general problem addressed by this research proposal is the design of non-supervised learning algorithms, in particular tensor factorization algorithms, applied in the space induced by a kernel function.

General objective

To design and implement an kernel-based tensor factorization algorithm for unsupervised learning.

Specific objectives

- To evaluate the impact, in terms of robustness and accuracy, of using particular kernels in kernel-based methods for tensor factorization.
- To design an efficient and effective kernel tensor factorization algorithms.
- To evaluate the performance of the proposed algorithms in both synthetic and real world datasets, and its impact in particular learning tasks: completion, clustering and blind source separation.
- To evaluate scalability comparing algorithms' complexity and parallelization perspectives.

6 Jusification

Due to the rich characteristics of natural processes and environments, it is rare that a single acquisition method provides complete understanding thereof. Information about a phenomenon or a system of interest can be obtained from different types of instruments, measurement techniques, experimental setups, and other types of sources. The increasing availability of multiple data sets that contain information, obtained using different acquisition methods, about the same system, introduces new degrees of freedom that raise questions beyond those related to analyzing each data set separately.

By treating that natural high-order array as a matrix, information is lost since lack the original multiway structure of the data.

Tensor factorizations have several advantages over two-way matrix factorizations including uniqueness of the optimal solution and component identification even when most of the data is missing. (Morup, 2015)

multiway decomposition techniques explicitly exploit the multiway structure that is lost when collapsing some of the modes of the tensor in order to analyze the data by regular matrix factorization approaches.

Tensor decompositions are in frequent use today in a variety of fields ranging from psychology, chemometrics, signal processing, bioinformatics, neuroscience, web mining, and computer vision to mention but a few.

Matrix and analogous tensor factorization is central to different important tasks in machine learning and information retrieval such as: clustering, latent topic analysis, recommendation, tensor completion, blind source separation, among others.

References

References