

# Data mining with tensor decompositions

Elaheh SOBHANI<sup>1</sup>, Pierre COMON<sup>1</sup>, Massoud BABAIE-ZADEH<sup>2</sup>

<sup>1</sup>GIPSA-Lab, Univ. Grenoble Alpes, CNRS, F-38000 Grenoble, France

<sup>2</sup>Sharif University of Technology, Tehran, Iran

prenom.nom@gipsa-lab.grenoble-inp.fr

**Résumé** – La fouille de textes est abordée en identifiant des modèles multi-vues. Suite aux travaux récents de A.Anandkumar, nous montrons que les algorithmes précédemment proposés présentent des inconvénients. L’algorithme que nous proposons permet de décomposer un tenseur non négatif construit avec des probabilités empiriques. Il se compare favorablement aux précédents sur des données synthétiques.

**Abstract** – Text mining is addressed by identifying multiview models. Following the recent works of A.Anandkumar, we show that the algorithms previously proposed exhibit important drawbacks. The algorithm subsequently proposed permits to decompose a nonnegative tensor built with sample joint probabilities. It compares favorably to the former on synthetically generated data sets.

## 1 Introduction

Multi-view models are useful in the context of multimodal recordings (such as text and audio), and offer a means to improve the estimation of latent variables that are present in several recordings. In this respect, they are related to data fusion. The model we assume does not impose any particular distribution (Gaussian, Dirichlet, etc). The notation is as follows.

Let  $L$  be the number of views,  $\mathbf{x}_\ell$  the observed views,  $1 \leq \ell \leq L$ , and  $h$  a discrete latent variable taking  $K$  possible integer values, say in  $\mathcal{H} = \{1, 2, \dots, K\}$  with probability  $\varphi(k) = \text{Prob}(h = k)$ . Every view belongs to a known dictionary  $\Omega = \{\mathbf{u}_1, \dots, \mathbf{u}_D\}$  of cardinality  $D$ , so that there exists a mapping  $\delta$  (not necessarily injective) from  $\mathcal{H}$  to  $\Omega$  such that  $\mathbf{x}_\ell = \mathbf{u}_{\delta(\ell)}$ . Then the key assumption is that views are statistically independent conditionally to  $h$ . More precisely, if we denote  $f_k(d) = \text{Prob}(\mathbf{x} = \mathbf{u}_d | h = k)$ , this means that the joint distribution of  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_L]$  can be written as:

$$p_{\mathbf{X}}(\mathbf{u}_{\delta(1)}, \dots, \mathbf{u}_{\delta(L)}) = \sum_{k=1}^K \varphi(k) f_k(\delta(1)) \dots f_k(\delta(L)) \quad (1)$$

The above equation is sometimes referred to as the *naive Bayes* model. The goal is to estimate the quantities appearing in the right-hand side of (1) from realizations of  $\mathbf{X}$ .

In the context of text mining,  $\mathbf{x}_\ell$  may correspond to words and  $h$  to a topic, for instance. As a function of the words appearing in a document containing  $L$  words, each belonging to a known dictionary of  $D$  possible words, it is thus possible to infer what is the probability that a topic is addressed therein.

The problem presentation we make in Sections 1-3 is the same as that described in [1,2], but adapted to the vocabulary of the Signal & Data Processing readership for the sake of clarity.

## 2 Generative Model

To generate a data set, we need first to define the distributions  $\varphi(k)$  and  $f_k(d)$  for all  $(k, \mathbf{u}_d) \in \mathcal{H} \times \Omega$ . The values of  $\varphi(k)$  are stored in a  $K$ -dimensional vector  $\boldsymbol{\varphi}$ . Similarly, the values of  $f_k(d)$  are stored in a  $D \times K$  matrix  $\mathbf{A}$ .

But it is actually useful to use their cumulative distributions,  $\Phi(k) = \text{Prob}(h \leq k)$  and  $F_k(d)$  with an appropriate encoding, as we shall see. The values of  $\Phi(k)$  are stored in a  $K \times 1$  vector  $\boldsymbol{\Phi}$  obtained from  $\boldsymbol{\varphi}$  by summing<sup>1</sup> its entries. Similarly, the values of  $F_k(d)$  are obtained by summing the entries of matrix  $\mathbf{A}$  and stored in a  $D \times K$  matrix  $\mathbf{F}$ . Our generative algorithm goes along the following lines:

- draw  $z \in [0, 1]$ , and pick  $h = \Phi^{-1}(z)$ , by selecting the first entry in  $\boldsymbol{\Phi}$  that is larger than  $z$ .
- for each  $\ell \in \{1, 2, \dots, L\}$ ,
  - draw  $z_\ell \in [0, 1]$ , and pick  $\delta(\ell) = F_k^{-1}(z_\ell)$ , by selecting the first entry in the  $k$ th column of  $\mathbf{F}$  that is larger than  $z_\ell$ .
  - set  $\mathbf{x}_\ell = \mathbf{u}_{\delta(\ell)}$ .

## 3 Moments

As in [1], to ease the presentation without restricting the generality, we may assume that  $\mathbf{u}_d$  are the columns of the  $D \times D$  identity matrix. Because of this choice made for  $\mathbf{u}_\ell$ , the joint probability (1) can be obtained by expectation as explained now. By definition of the encoding of  $\mathbf{x}_\ell$ , we have:

$$\mathbf{x}_p \otimes \mathbf{x}_q \otimes \mathbf{x}_r = \mathbf{u}_{\delta(p)} \otimes \mathbf{u}_{\delta(q)} \otimes \mathbf{u}_{\delta(r)} \quad (2)$$

This tensor has only one nonzero entry, which means that by summing  $N$  realizations, we obtain an estimate of the number

<sup>1</sup>With this goal, the matlab function `cumsum()` can be used.

of occurrences. Hence averaging  $N$  realizations yields the joint probability distribution. In other words,

$$\hat{p}_{x_p, x_q, x_r}(\mathbf{u}_{\delta(p)}, \mathbf{u}_{\delta(q)}, \mathbf{u}_{\delta(r)}) = \widehat{\mathbb{E}}\{\mathbf{x}_p \otimes \mathbf{x}_q \otimes \mathbf{x}_r\} \Big|_{\delta(p), \delta(q), \delta(r)} \quad (3)$$

where  $(\hat{\cdot})$  denotes sample estimate. In the sequel, we shall need the moments of order two and three, which correspond to double and triple joint probabilities, denoted as follows:

$$\mathbf{P} = \mathbb{E}\{\mathbf{x}_p \otimes \mathbf{x}_q\} \quad (4)$$

$$\mathcal{T} = \mathbb{E}\{\mathbf{x}_p \otimes \mathbf{x}_q \otimes \mathbf{x}_r\} \quad (5)$$

where  $\mathbf{P}$  is a  $D \times D$  symmetric matrix and  $\mathcal{T}$  a  $D \times D \times D$  symmetric tensor. These moments do not depend on  $\{p, q, r\}$  provided these 3 integers are all different, which ensures the conditional independence assumed in (1). Note that  $\{\delta(p), \delta(q), \delta(r)\}$  may not be different because  $\delta$  is not injective.

## 4 Previous approaches

We describe in this section the approaches proposed by A.Anandkumar [1, 2]. These approaches may seem attractive at first glance, but it turns out that they are not usable in practice for the reasons we subsequently point out. In fact, there is no computer experiment in the latter papers, which may explain why the serious problems encountered in these algorithms have not been detected.

### 4.1 Power method with deflation

In [1], the authors propose to use the two moments defined in (4) and (5). Because of the Bayesian rule (1) and (3), these moments enjoy the following relations:

$$\mathbf{P} = \sum_{k=1}^K \varphi_k \mathbf{a}_k \otimes \mathbf{a}_k \quad (6)$$

$$\mathcal{T} = \sum_{k=1}^K \varphi_k \mathbf{a}_k \otimes \mathbf{a}_k \otimes \mathbf{a}_k \quad (7)$$

where  $\mathbf{a}_k$  denotes the  $k$ th column of matrix  $\mathbf{A}$  defined in Section 2. Because of (6), matrix  $\mathbf{P}$  is theoretically positive semidefinite, since  $\varphi_k$  are positive numbers. Hence, in a similar way as had been done for Blind Source Separation [3], there exists a “whitening” matrix  $\mathbf{W}$  such that  $\mathbf{W}^T \mathbf{P} \mathbf{W} = \mathbf{I}$ , where  $\mathbf{I}$  is the identity matrix;  $\mathbf{W}$  can theoretically be easily obtained from the EigenValue Decomposition (EVD)  $\mathbf{P} = \mathbf{U} \mathbf{D} \mathbf{U}^T$ ,  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ , by setting  $\mathbf{W} = \mathbf{U} \mathbf{D}^{-1/2}$ . Then we have

$$\sum_{k=1}^K \tilde{\mathbf{a}}_k \tilde{\mathbf{a}}_k^T = \mathbf{I}$$

if  $\tilde{\mathbf{a}}_k \stackrel{\text{def}}{=} \sqrt{\varphi_k} \mathbf{W}^T \mathbf{a}_k$ . In other words, the matrix  $\tilde{\mathbf{A}}$  containing  $\tilde{\mathbf{a}}_k$  as columns is now  $K \times K$  orthogonal:  $\tilde{\mathbf{A}} \tilde{\mathbf{A}}^T = \mathbf{I}$ . However, if matrix  $\mathbf{P}$  is estimated via the mere averaging (4) as proposed in [1],  $\mathbf{P}$  may have negative eigenvalues because of estimation errors; this has been overlooked in [1].

Next, this whitening matrix is applied to tensor  $\mathcal{T}$  to yield  $\tilde{\mathcal{T}} \stackrel{\text{def}}{=} \mathcal{T} \bullet \mathbf{W} \bullet \mathbf{W} \bullet \mathbf{W}$ . This new tensor enjoys the relationship

$$\tilde{\mathcal{T}} = \sum_{k=1}^K \varphi_k^{-1/2} \tilde{\mathbf{a}}_k \otimes \tilde{\mathbf{a}}_k \otimes \tilde{\mathbf{a}}_k$$

The conclusion is that  $\tilde{\mathcal{T}}$  ideally admits an *orthogonal* Canonical Polyadic (CP) decomposition; see *e.g.* [4] for an introduction. Many algorithms have been devised for this task, including the pair-sweeping CoM algorithm, or Joint Approximate Diagonalization (JAD) algorithms [3]. But it is proposed in [1] to execute the tensor power iteration [5] to extract the dominant “eigenvector<sup>2</sup>” of  $\tilde{\mathcal{T}}$  and then proceed by deflation to get the remaining ones.

One known problem with deflation, which has been observed even for matrices, is that the eigenvectors extracted may loose orthogonality because of rounding errors. The consequence is that the same (dominant) eigenvector may show up several times. This has been also overlooked in [1]. We propose in Section 5 to fix the problem by a re-orthogonalization with previously found vectors. For small size problem, drawing an initial value in the orthogonal subspace is sufficient.

### 4.2 Diagonalization of two moment matrices

In [2], a joint diagonalization algorithm is promoted and uses two moment matrices, namely  $\mathbf{P}$  and a matrix contraction  $\mathbf{T}(\boldsymbol{\eta}) = \mathcal{T} \bullet \boldsymbol{\eta}$ , where  $\boldsymbol{\eta}$  is a randomly drawn vector. The idea is very interesting, but the algorithm unfortunately has never been implemented and tested, according to the available literature. We describe it shortly in this section, and shall not report its poor performances in Section 6 for reasons of space.

The Algorithm A of [2] starts with an SVD of  $\mathbf{P}$  as  $\mathbf{P} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T$ . Then the matrix  $\mathbf{B}(\boldsymbol{\eta}) = \mathbf{U}^T \mathbf{T}(\boldsymbol{\eta}) \mathbf{V} (\mathbf{U}^T \mathbf{P} \mathbf{V})^{-1}$  is computed, as well as its  $K$  dominant eigenvectors,  $\boldsymbol{\xi}_k$ . Then an estimate of the columns of matrix  $\mathbf{A}$  is given by  $\hat{\mathbf{a}}_k = \mathbf{U} \boldsymbol{\xi}_k$ , up to a scaling factor depending on how  $\hat{\mathbf{a}}_k$  are normalized. The “eigenvalues” of  $\mathcal{T}$  can be obtained in a second stage by contraction as  $\hat{\varphi}_k = \mathcal{T} \bullet \hat{\mathbf{a}}_k \bullet \hat{\mathbf{a}}_k \bullet \hat{\mathbf{a}}_k$ .

## 5 The algorithms proposed

### 5.1 Re-orthogonalization of the Power Method

Let  $\mathbf{V}(i) \stackrel{\text{def}}{=} [\tilde{\mathbf{a}}_1, \dots, \tilde{\mathbf{a}}_i]$  be the matrix containing the eigenvectors already output at iteration  $i$ . To compute the next iteration, we draw a random initial vector  $\boldsymbol{\theta}$  and orthogonalize it with respect to  $\mathbf{V}$  by assuming  $\tilde{\boldsymbol{\theta}} = \boldsymbol{\theta} - \mathbf{V} \mathbf{V}^T \boldsymbol{\theta}$  as initial vector, instead of  $\boldsymbol{\theta}$  as in [1]. This modification is sufficient in small dimensions. However, in larger dimensions, it will be necessary

<sup>2</sup>Recall that there exist several definitions of tensor eigenvectors. The definition used in [1] – and hence here – is  $\mathcal{T} \bullet \mathbf{v} \bullet \mathbf{v} = \lambda \mathbf{v}$ , which has the undesirable property that  $\lambda$  depends on the norm of  $\mathbf{v}$ . In fact, if  $(\lambda, \mathbf{v})$  is an eigenpair, then so is  $(\alpha \lambda, \alpha \mathbf{v})$  for any nonzero  $\alpha$ .

to also re-orthogonalize the intermediate iterates generated by the Power Method because of rounding errors.

## 5.2 Algorithm under nonnegative constraint

As can be seen in (7),  $\mathcal{T}$  has a CP decomposition structure. Since we expect to obtain  $\mathbf{A}$  and  $\varphi$ , which contain probability values, nonnegativity is an essential constraint that should be taken into account during the decomposition. Although (7) reveals that  $\mathcal{T}$  is a symmetric tensor, we do not need to impose symmetry during the course of the CP decomposition.

CP decomposition viewed as a sort of low-rank approximation is ill-posed in general, unless some constraints such as nonnegativity are imposed [4]. Therefore, we consider the following minimization in order to obtain reliable estimated quantities,  $\mathbf{A}$  and  $\varphi$ , in  $\mathbb{R}_+$ :

$$\begin{aligned} \min_{\mathbf{A}, \varphi} & \|\mathcal{T} - \sum_{k=1}^K \varphi(k) \mathbf{a}_k \otimes \mathbf{a}_k \otimes \mathbf{a}_k\|_2 \\ \text{s.t.} & \mathbf{A} \in \mathbb{R}_+^{D \times K}, \varphi \in \mathbb{R}_+^K. \end{aligned} \quad (8)$$

Obviously, (8) is a nonconvex problem, and Alternating Optimization (AO) is usually utilized for such a problem [6–8]. AO is also known as Block Coordinate Descent (BCD) and its convergence has been discussed in [9]. Alternating Least Squares (ALS) is a special case of AO when the cost function is simply least squares with no regularization [6].

---

**Algorithm 1** AO-ADMM for Problem (8)

---

**Input:**  $\mathcal{T}, K, \rho$

**Output:**  $\mathbf{A}, \varphi$

- 1: Initialize  $\mathbf{H}_1, \mathbf{H}_2, \mathbf{H}_3$
  - 2: Initialize  $\mathbf{U}_1, \mathbf{U}_2, \mathbf{U}_3$
  - 3: **repeat**
  - 4:   **for**  $n = 1, 2, 3$  **do**
  - 5:      $\mathbf{T} = \mathcal{T}_{(n)}$  (Unfold mode  $n$ )
  - 6:      $\mathbf{Z} = \odot_{j \neq n} \mathbf{H}_j$
  - 7:     Compute the Cholesky factor  $\mathbf{L}$  of  $(\mathbf{Z}^\top \mathbf{Z} + \rho \mathbf{I})$
  - 8:     Update  $\mathbf{H}_n$  as in Alg.1 of [6]:
    - (a)  $\mathbf{G}_n \leftarrow \mathbf{L}^{-\top} \mathbf{L}^{-1} [\mathbf{Z}^\top \mathbf{T} + \rho(\mathbf{U}_n + \mathbf{H}_n)^\top]$  # Auxiliary variables
    - (b)  $\mathbf{H}_n \leftarrow \max[\mathbf{0}, \mathbf{G}_n^\top - \mathbf{U}_n]$  # Primal variables
    - (c)  $\mathbf{U}_n \leftarrow \mathbf{U}_n + \mathbf{H}_n - \mathbf{G}_n^\top$  # Dual variables
  - 9:   **end for**
  - 10: **until** some termination criterion
  - 11: **for**  $k = 1 : K$  **do**
  - 12:    $\mathbf{a}_k = \mathbf{H}_1(:, k) / \|\mathbf{H}_1(:, k)\|_1; \varphi_k = \prod_n \|\mathbf{H}_n(:, k)\|_1$
  - 13: **end for**
  - 14:  $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_K], \varphi = [\varphi_1, \varphi_2, \dots, \varphi_K]^\top$
- 

Recently, a new method has been introduced in [6], which improves AO by the means of Alternating Direction Method of Multipliers (ADMM) [10], and can handle constraints. This

new method is called AO-ADMM [6]. Applying AO-ADMM to our problem is straightforward, and its details are shown in Algorithm 1. The only (mild) limitation of Algorithm 1 is related to the uniqueness of the CP decomposition. From (7), the rank of  $\mathcal{T}$  is  $K$ , which should be less than the *expected rank*  $R^o$  to ensure uniqueness, or than  $R_s^o$  in the symmetric case [4], with:

$$R^o = \left\lceil \frac{D^3}{3D-2} \right\rceil, R_s^o = \left\lceil \frac{(D+1)(D+2)}{6} \right\rceil.$$

Hence, we have the condition,  $K = O(D^2)$  for the proposed algorithm. All in all, we only need moments of order three, *i.e.*  $\mathcal{T}$ , to obtain probabilities of hidden multi-view variables, and addressing over-complete regimes up to  $K = O(D^2)$  is possible (contrary to the Power Method).

## 6 Computer experiments

We explained in Section 2 how one can generate a set of multi-view variables that are related to the same hidden variable, with these two properties simultaneously:

- conditionally independent given the hidden variable
- having the same conditional distribution

In order to obtain an acceptable approximation of second and third order moments, we assume a large number of trials, say  $N=100,000$ . In each trial, we draw a random hidden variable,  $h$ , in the way described in Section 2. Then, for the chosen hidden variable, we draw three random variable,  $(x_p, x_q, x_r)$ , such that they satisfy the two properties mentioned above (cf. Sections 1-2). At the end, by averaging arrays  $\mathbf{P}$  and  $\mathcal{T}$  obtained in each trial, we have an approximation of second and third order moments respectively.

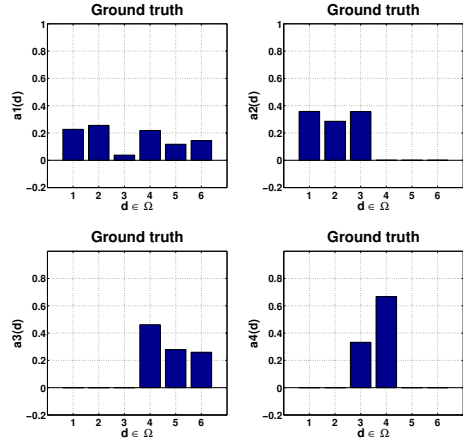


FIG. 1: Original  $\mathbf{A}$

For the sake of representation, we consider  $K = 4$ , thereby we generate four distinct distributions which are actually the columns of  $\mathbf{A}$ . Because of the lack of space, we able to show only the comparison between various methods with only one single size of dictionary,  $D = 6$ , so that  $\mathbf{A}$  is  $6 \times 4$ .

Figure 1 (Left) shows the values of the  $D = 6$  conditional probabilities contained in each  $\mathbf{a}_k$ ,  $1 \leq k \leq K$ . Fig. 1 (Right)

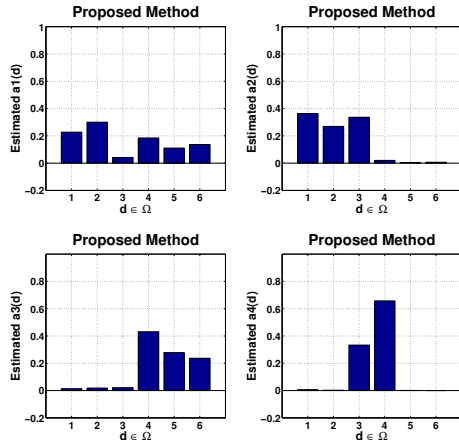


FIG. 2: Method proposed in Section 5.2

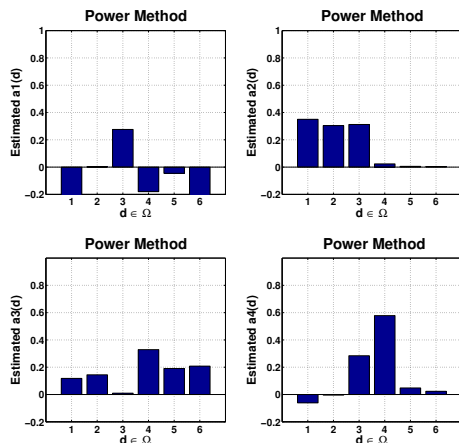


FIG. 3: The Robust tensor power method [1]

reports the corresponding estimated values obtained by the algorithm described in Section 5.2 (after optimal permutation to get them in the same order). As shown in Fig. 4 (Left), the Power method fails in delivering the correct columns of  $\mathbf{A}$ , and suffers from negative values. This is the main difference between [1] and the proposed algorithm which imposes the non-negative constraint. The performances in estimating vector  $\varphi$  have been compared in Fig. 4 (Right). It is clear that the result of proposed method is closer to the original vector  $\varphi$ .

## 7 Discussion

Estimates of  $\mathbf{P}$  and  $\mathcal{T}$  from data via (4) and (5) may not yield symmetric arrays, which raises serious problems. It can be fixed by forcing the symmetry of the estimate of  $\mathbf{P}$  before running the Power Method. Next, the Power Method is not the most adequate algorithm to compute a CP decomposition, even under the orthogonal constraints; in fact, it is very sensitive to initialization<sup>3</sup>, and may output several times the same dominant eigenvector. This can be faced by the re-orthogonalization pro-

<sup>3</sup>In [1], many initial values are tried, and the best one is picked. In Section 6, the Power Method has been initialized with “only” 10 initial values for each eigenvector, to keep a reasonable computational complexity.

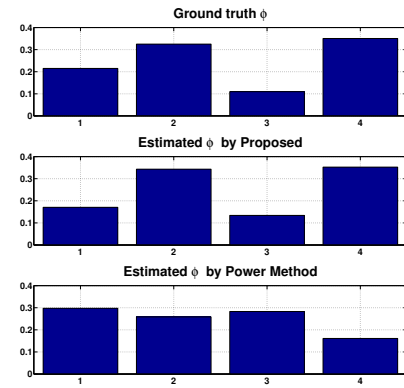


FIG. 4: Performance of methods in estimating  $\varphi$

cedure described in Section 5.1. A more serious limitation of the algorithms described in Section 4 is the output of negative values. The algorithm proposed in Section 5 fixes the problem, as demonstrated in Section 6. In addition, it can handle over-complete cases, where  $K > D$ , provided  $K$  is not larger than  $O(D^2)$ , whereas the Power Method cannot.

## References

- [1] A. Anandkumar, R. Ge, D. Hsu, S. M. Kakade, and M. Telgarsky, “Tensor decompositions for learning latent variable models,” *J. Machine Learning Research*, vol. 15, pp. 2773–2832, Aug. 2014.
- [2] A. Anandkumar, D. Hsu, and S. M. Kakade, “A method of moments for mixture models and hidden markov models,” in *25th Annual Conference on learning Theory*, vol. 23, 2012, pp. 33.1–33.34. [Online]. Available: arxiv:1203.0683
- [3] P. Comon and C. Jutten, Eds., *Handbook of Blind Source Separation*. Oxford UK, Burlington USA: Academic Press, 2010.
- [4] P. Comon, “Tensors: a brief introduction,” *IEEE Signal Processing Magazine*, vol. 31, no. 3, pp. 44–53, 2014.
- [5] L. DeLathauwer, P. Comon *et al.*, “Higher-order power method, application in Independent Component Analysis,” in *NOLTA*, Las Vegas, Dec. 1995, pp. 91–96.
- [6] K. Huang, N. Sidiropoulos, and A. P. Liavas, “A flexible and efficient algorithmic framework for constrained matrix and tensor factorization,” *IEEE Trans. Signal proc.*, vol. 64, no. 19, pp. 5052–5065, 2016.
- [7] E. Sobhani, M. Sadeghi, M. Babaie-Zadeh, and C. Jutten, “A robust ellipse fitting algorithm based on sparsity of outliers,” in *2017 25th European Signal Processing Conference (EUSIPCO)*. IEEE, 2017, pp. 1195–1199.
- [8] P. Comon, X. Luciani, and A. L. De Almeida, “Tensor decompositions, alternating least squares and other tales,” *Jour. Chemometrics*, vol. 23, no. 7-8, pp. 393–405, 2009.
- [9] A. Beck and L. Tetrushvili, “On the convergence of block coordinate descent type methods,” *SIAM journal on Optimization*, vol. 23, no. 4, pp. 2037–2060, 2013.
- [10] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, “Distributed optimization and statistical learning via the alternating direction method of multipliers,” *Foundations and Trends in Machine learning*, vol. 3, no. 1, pp. 1–122, 2011.