

NON-NEGATIVE TENSOR FACTORIZATION USING ALPHA AND BETA DIVERGENCES

Andrzej CICHOCKI^{1*}, Rafal ZDUNEK^{1†}, Seungjin CHOI², Robert PLEMMONS³ and Shun-ichi AMARI¹

¹ Brain Science Institute, RIKEN, Wako-shi, Saitama 351-0198, JAPAN,

² Pohang University of Science and Technology, KOREA,

³ Wake Forest University, USA

ABSTRACT

In this paper we propose new algorithms for 3D tensor decomposition/factorization with many potential applications, especially in multi-way Blind Source Separation (BSS), multidimensional data analysis, and sparse signal/image representations. We derive, compare and implement in MATLAB NTFLAB Toolbox three classes of algorithms: Multiplicative, Fixed Point Alternating Least Squares (FPALS) and Alternating Interior-Point Gradient (AIPG) algorithms. Some of the proposed algorithms are characterized by improved robustness, efficiency and convergence rates and can be applied for various distributions of data and additive noise.

Index Terms— Algorithms, Learning systems, Linear approximation, Signal representations, Feature extraction.

1. MODELS AND PROBLEM FORMULATION

Tensors (also known as n-way arrays or multidimensional arrays) are used in a variety of applications ranging from neuroscience and psychometrics to chemometrics [6,8,9,17-19]. Nonnegative matrix factorization (NMF), Non-negative tensor factorization (NTF) and parallel factor analysis PARAFAC models with non-negativity constraints have been recently proposed as promising sparse and quite efficient representations of signals, images, or general data [2-7,10-13]. From a viewpoint of data analysis, NTF is very attractive because it takes into account spacial and temporal correlations between variables more accurately than 2D matrix factorizations, such as NMF, and it provides usually sparse common factors or hidden (latent) components with physiological meaning and interpretation [9,15]. In most applications, especially in neuroscience (EEG, fMRI), the standard NTF or PARAFAC models were used [15,16]. In this paper we consider more general model referred to as 3D NTF2 model (in analogy to the Parafac2 model [17]) (see Fig. 1). A given tensor $\underline{\mathbf{X}} \in \mathbb{R}_+^{I \times T \times K}$ is decomposed to a set of matrices \mathbf{S} , \mathbf{D} and $\{\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_K\}$ with nonnegative entries. Here and elsewhere, \mathbb{R}_+ denotes the nonnegative orthant with appropriate

dimensions. The three-way NTF2 model can be described as

$$\mathbf{X}_k = \mathbf{A}_k \mathbf{D}_k \mathbf{S} + \mathbf{E}_k, \quad (k = 1, 2, \dots, K) \quad (1)$$

where $\mathbf{X}_k = \mathbf{X}_{::,k} = [x_{itk}]_{I \times T} \in \mathbb{R}_+^{I \times T}$ are frontal slices of $\underline{\mathbf{X}} \in \mathbb{R}_+^{I \times T \times K}$, K is the number of frontal slices, $\mathbf{A}_k = [a_{irk}]_{I \times R} \in \mathbb{R}_+^{I \times R}$ are the basis (mixing matrices), $\mathbf{D}_k \in \mathbb{R}_+^{R \times R}$ is a diagonal matrix that holds the k -th row of the $\mathbf{D} \in \mathbb{R}_+^{K \times R}$ in its main diagonal, and $\mathbf{S} = [s_{rt}]_{R \times T} \in \mathbb{R}_+^{R \times T}$ is a matrix representing sources (or hidden components or common factors), and $\mathbf{E}_k = \mathbf{E}_{::,k} \in \mathbb{R}_+^{I \times T}$ is the k -th frontal slice of a tensor $\underline{\mathbf{E}} \in \mathbb{R}_+^{I \times T \times K}$ representing error or noise depending upon the application. The objective is to estimate the set of matrices $\{\mathbf{A}_k\}$, (k, \dots, K) , \mathbf{D} and \mathbf{S} , subject to some non-negativity constraints and other possible natural constraints such as sparseness and/or smoothness. Since the diagonal matrices \mathbf{D}_k are scaling matrices they can usually be absorbed by the matrices \mathbf{A}_k by introducing column-normalized matrices $\mathbf{A}_k := \mathbf{A}_k \mathbf{D}_k$, so usually in BSS applications the matrix \mathbf{S} and the set of scaled matrices $\mathbf{A}_1, \dots, \mathbf{A}_K$ need only to be estimated. However, in such a case we may lose the uniqueness of the NTF representation ignoring scaling and permutation ambiguities. The uniqueness still can be achieved by imposing nonnegativity, sparsity and other constraints. The above NTF2 model is similar to the well known PARAFAC2 model with non-negativity constraints and Tucker models [6,15,17]. In the special case, when all matrices \mathbf{A}_k are identical, the NTF2 model can be simplified to the ordinary PARAFAC model with the non-negativity constraints described $\mathbf{X}_k = \mathbf{A} \mathbf{D}_k \mathbf{S} + \mathbf{E}_k$, $k = 1, \dots, K$ or equivalently $x_{itk} = \sum_r a_{ir} s_{rt} d_{kr} + e_{itk}$ or $\underline{\mathbf{X}} = \sum_r \mathbf{a}_r \otimes \mathbf{s}_r^T \otimes \mathbf{d}_r + \underline{\mathbf{E}}$, where \mathbf{s}_r are rows of \mathbf{S} and $\mathbf{a}_r, \mathbf{d}_r$ are columns of \mathbf{A} and \mathbf{D} , respectively and \otimes means outer product of vectors [8,9]. Throughout this paper, we use the following notation: the rt -th element of the matrix \mathbf{S} is denoted by s_{rt} , $x_{itk} = [\mathbf{X}_k]_{it}$ means the it -th element of the k -th frontal slice matrix \mathbf{X}_k , $\bar{\mathbf{A}} = [\mathbf{A}_1; \mathbf{A}_2; \dots; \mathbf{A}_K] \in \mathbb{R}_+^{K \times I \times R}$ is a column-wise unfolded matrix of the slices \mathbf{A}_k , $\bar{a}_{pr} = [\bar{\mathbf{A}}]_{pr}$. Analogously, $\bar{\mathbf{X}} = [\mathbf{X}_1; \mathbf{X}_2; \dots; \mathbf{X}_K] \in \mathbb{R}_+^{K \times I \times T}$ is the column-wise unfolded matrix of the slices \mathbf{X}_k and $\bar{x}_{pt} = [\bar{\mathbf{X}}]_{pt}$.

* On leave from Warsaw University of Technology, Poland

† On the leave from Institute of Telecommunications, Teleinformatics, and Acoustics, Wrocław University of Technology, Poland

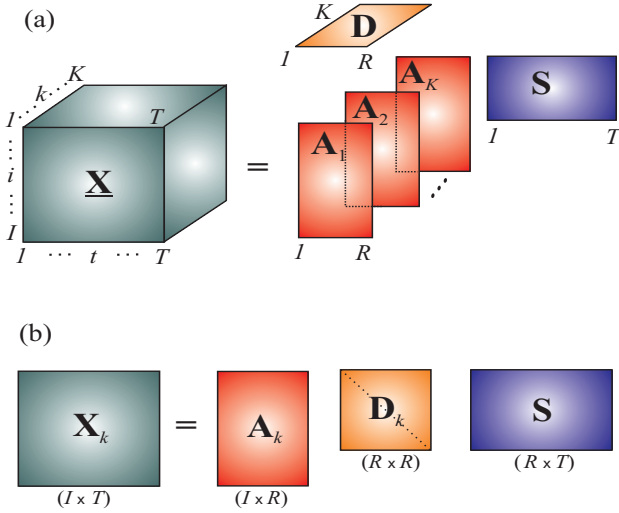


Fig. 1. (a) NTF2 model in which 3D tensor is decomposed to set of nonnegative matrices: $\{\mathbf{A}_1, \dots, \mathbf{A}_K\}$, \mathbf{D} , \mathbf{S} . (b) Equivalent representation in which frontal slices of tensor are factored by set of matrices (tensor $\underline{\mathbf{E}}$ representing error is omitted for simplicity).

2. ALPHA AND BETA DIVERGENCES AND ASSOCIATED ALGORITHMS

To deal with the model (1) efficiently we adopt several approaches from constrained optimization and multi-criteria optimization, where we minimize simultaneously several cost functions using alternating switching between sets of parameters. Alpha and Beta divergences are two complimentary generalized cost functions which can be applied for NMF and NTF [1,3,4,5,20].

2.1. Alpha Divergence

Let us consider a flexible and general class of the cost functions, called α -divergence [1,3,4]:

$$D_A^{(\alpha)}(\bar{\mathbf{X}} \parallel \bar{\mathbf{A}}\mathbf{S}) = \frac{\sum_{pt} (\bar{x}_{pt}^\alpha [\bar{\mathbf{A}}\mathbf{S}]_{pt}^{1-\alpha} - \alpha \bar{x}_{pt} + (\alpha - 1) [\bar{\mathbf{A}}\mathbf{S}]_{pt})}{\alpha(\alpha - 1)}$$

$$\begin{aligned} D_{A_k}^{(\alpha)}(\mathbf{X}_k \parallel \mathbf{A}_k \mathbf{S}) &= \\ &= \frac{\sum_{it} (x_{itk}^\alpha [\mathbf{A}_k \mathbf{S}]_{it}^{1-\alpha} - \alpha x_{itk} + (\alpha - 1) [\mathbf{A}_k \mathbf{S}]_{it})}{\alpha(\alpha - 1)} \end{aligned}$$

We note that as special cases of α -divergence for $\alpha = 2, 0.5, -1$, we obtain the Pearson's, Hellinger's and Neyman's chi-square distances, respectively, while for the cases $\alpha = 1$ and $\alpha = 0$ the divergence has to be defined by the limits: $\alpha \rightarrow 1$ and $\alpha \rightarrow 0$, respectively. When these limits are evaluated one obtains for $\alpha \rightarrow 1$ the generalized Kullback-Leibler divergence (I-divergence) and for $\alpha \rightarrow 0$ the dual generalized KL divergence [1,3,4].

Instead of applying the standard gradient descent method, we use the nonlinearly transformed gradient approach as generalization of the exponentiated gradient (EG)[4]:

$$\Phi(a_{irk}) \leftarrow \Phi(a_{irk}) - \eta_{irk} \frac{\partial D_{A_k}^{(\alpha)}(\mathbf{X}_k \parallel \mathbf{A}_k \mathbf{S})}{\partial \Phi(a_{irk})}, \quad (2)$$

$$\Phi(s_{rt}) \leftarrow \Phi(s_{rt}) - \eta_{rt} \frac{\partial D_A^{(\alpha)}(\bar{\mathbf{X}} \parallel \bar{\mathbf{A}}\mathbf{S})}{\partial \Phi(s_{rt})}, \quad (3)$$

where $\Phi(x)$ is a suitably chosen function.

It can be shown that such a nonlinear scaling or transformation provides a stable solution and the gradients are much better behaved in the space Φ . In our case, we employ $\Phi(x) = x^\alpha$, which leads directly to the new learning algorithm (for $\alpha \neq 0$) (the rigorous proof of local convergence similar to this given by Lee and Seung [12] is omitted due to a lack of space):

$$a_{irk} \leftarrow a_{irk} \left(\frac{\sum_{t=1}^T (x_{itk} / [\mathbf{A}_k \mathbf{S}]_{it})^\alpha s_{rt}}{\sum_{t=1}^T s_{rt}} \right)^{1/\alpha}, \quad (4)$$

$$s_{rt} \leftarrow s_{rt} \left(\frac{\sum_{p=1}^{KI} \bar{a}_{pr} (\bar{x}_{pt} / [\bar{\mathbf{A}}\mathbf{S}]_{pt})^\alpha}{\sum_{p=1}^{KI} \bar{a}_{pr}} \right)^{1/\alpha}. \quad (5)$$

2.2. Beta Divergence

Regularized beta divergence for the NTF2 problem can be defined as follows:

$$\begin{aligned} D^{(\beta)}(\bar{\mathbf{X}} \parallel \bar{\mathbf{A}}\mathbf{S}) &= \sum_{pt} (\bar{x}_{pt} \frac{\bar{x}_{pt}^\beta - [\bar{\mathbf{A}}\mathbf{S}]_{pt}^\beta}{\beta(\beta + 1)} \\ &+ [\bar{\mathbf{A}}\mathbf{S}]_{pt}^\beta \frac{[\bar{\mathbf{A}}\mathbf{S}]_{pt} - \bar{x}_{pt}}{\beta + 1}) + \alpha_S \|\mathbf{S}\|_{L1}, \end{aligned} \quad (6)$$

$$\begin{aligned} D_k^{(\beta)}(\mathbf{X}_k \parallel \mathbf{A}_k \mathbf{S}) &= \sum_{it} (x_{itk} \frac{x_{itk}^\beta - [\mathbf{A}_k \mathbf{S}]_{it}^\beta}{\beta(\beta + 1)} \\ &+ [\mathbf{A}_k \mathbf{S}]_{it}^\beta \frac{[\mathbf{A}_k \mathbf{S}]_{it} - x_{itk}}{\beta + 1}) + \alpha_{A_k} \|\mathbf{A}_k\|_{L1}, \end{aligned} \quad (7)$$

$k = 1, \dots, K, \quad t = 1, 2, \dots, T, \quad i = 1, 2, \dots, I,$

where α_S and α_{A_k} are small positive regularization parameters which control the degree of sparseness of the matrices \mathbf{S} and \mathbf{A}_k , respectively, and the $L1$ -norms defined as $\|\mathbf{S}\|_{L1} = \sum_{rt} |s_{rt}|$ and $\|\mathbf{A}_k\|_{L1} = \sum_{ir} |a_{irk}|$ are introduced to enforce sparse representations of the solutions. It is interesting to note that for $\beta = 1$, we obtain the squared Euclidean distances expressed by the Frobenius norms $\|\mathbf{X}_k - \mathbf{A}_k \mathbf{S}\|_F^2$, while for the singular cases, $\beta = 0$ and $\beta = -1$, the beta divergence has to be defined as limiting cases as $\beta \rightarrow 0$ and $\beta \rightarrow -1$, respectively. When these limits are evaluated one gets for $\beta \rightarrow 0$ the generalized Kullback-Leibler divergence (called I-divergence) and for $\beta \rightarrow -1$ we obtain the Itakura-Saito distance. The choice of the parameter β depends on the statistical distribution of the data and the beta

divergence corresponds to the Tweedie models. For example, the optimal choice of the parameter for the normal distribution is $\beta = 1$, for the gamma distribution is $\beta \rightarrow -1$, for the Poisson distribution $\beta \rightarrow 0$, and for the compound Poisson $\beta \in (-1, 0)$. By minimizing the beta divergence, we have derived various kinds of NTF algorithms: Multiplicative based on the standard gradient descent, Exponentiated Gradient (EG), Projected Gradient (PG), Alternating Interior-Point Gradient (AIPG), or Fixed Point Alternating Least Squares (FPALS) algorithms. For example, in order to derive a flexible multiplicative learning algorithm, we compute the gradient of (6)-(7) with respect to elements of matrices $s_{rt} = s_r(t) = [\mathbf{S}]_{rt}$ and $a_{irk} = [\mathbf{A}_k]_{ir}$ and performing simple mathematical manipulations we obtain the multiplicative update rules:

$$a_{irk} \leftarrow a_{irk} \frac{[\sum_{t=1}^T (x_{itk} / [\mathbf{A}_k \mathbf{S}]_{it}^{1-\beta}) s_{rt} - \alpha_{A_k}]_{\varepsilon}}{\sum_{t=1}^T [\mathbf{A}_k \mathbf{S}]_{it}^{\beta} s_{rt}}, \quad (8)$$

$$s_{rt} \leftarrow s_{rt} \frac{[\sum_{p=1}^{KI} \bar{a}_{pr} (\bar{x}_{itk} / [\bar{\mathbf{A}} \mathbf{S}]_{it}^{1-\beta}) - \alpha_S]_{\varepsilon}}{\sum_{p=1}^{KI} \bar{a}_{pr} [\bar{\mathbf{A}} \mathbf{S}]_{it}^{\beta}}, \quad (9)$$

where $[x]_{\varepsilon} = \max\{\varepsilon, x\}$ with a small $\varepsilon = 10^{-16}$ is introduced in order to avoid zero and negative values.

In the special case for $\beta = 1$ we can derive an alternative algorithm, called FPALS (Fixed Point Alternating Least Squares) algorithm (see [5] for detail)

$$\mathbf{A}_k \leftarrow \left[(\mathbf{X}_k \mathbf{S}^T - \alpha_{A_k} \mathbf{E}_A) (\mathbf{S} \mathbf{S}^T)^+ \right]_{\varepsilon}, \quad (10)$$

$$\mathbf{S} \leftarrow \left[(\bar{\mathbf{A}}^T \bar{\mathbf{A}})^+ (\bar{\mathbf{A}}^T \bar{\mathbf{X}} - \alpha_S \mathbf{E}_S) \right]_{\varepsilon}, \quad (11)$$

where $[\mathbf{A}]^+$ denotes Moore-Penrose pseudo-inverse of \mathbf{A} and $\mathbf{E}_A \in \mathbb{R}^{I \times R}$, $\mathbf{E}_S \in \mathbb{R}^{R \times T}$ are matrices with all entries one. The above algorithm can be considered as a nonlinear projected Alternating Least Squares (ALS) or nonlinear extension of EM-PCA algorithm.

Furthermore, using the Alternating Interior-Point Gradient (AIPG) approach [14], another efficient algorithm has been developed and implemented [5]:

$$\mathbf{A}_k \leftarrow \mathbf{A}_k - \eta_{A_k} \mathbf{P}_{A_k}, \quad (12)$$

$$\mathbf{S} \leftarrow \mathbf{S} - \eta_S \mathbf{P}_S, \quad (13)$$

where $\mathbf{P}_{A_k} = (\mathbf{A}_k \odot (\mathbf{A}_k \mathbf{S} \mathbf{S}^T)) \oslash ((\mathbf{A}_k \mathbf{S} - \mathbf{X}_k) \mathbf{S}^T)$, $\mathbf{P}_S = (\mathbf{S} \odot (\bar{\mathbf{A}}^T \bar{\mathbf{A}} \mathbf{S})) \oslash (\bar{\mathbf{A}}^T (\bar{\mathbf{A}} \mathbf{S} - \bar{\mathbf{X}}))$ and operators \odot and \oslash mean component-wise multiplication and division, respectively. The learning rates η_{A_k} and η_S are selected in this way to ensure the steepest descent, and on the other hand, to maintain non-negativity. Thus, $\eta_S = \min\{\tau \hat{\eta}_S, \eta_S^*\}$ and $\eta_{A_k} = \min\{\tau \hat{\eta}_{A_k}, \eta_{A_k}^*\}$, where $\tau \in (0, 1)$, $\hat{\eta}_S = \{\eta : \mathbf{S} - \eta \mathbf{P}_S\}$ and $\hat{\eta}_{A_k} = \{\eta : \mathbf{A}_k - \eta \mathbf{P}_{A_k}\}$ ensure non-negativity,

and

$$\eta_{A_k}^* = \frac{\text{vec}(\mathbf{P}_{A_k})^T \text{vec}(\mathbf{A}_k \mathbf{S} \mathbf{S}^T - \mathbf{X}_k \mathbf{S}^T)}{\text{vec}(\mathbf{P}_{A_k} \mathbf{S})^T \text{vec}(\mathbf{P}_{A_k} \mathbf{S})}, \quad (14)$$

$$\eta_S^* = \frac{\text{vec}(\mathbf{P}_S)^T \text{vec}(\bar{\mathbf{A}}^T \bar{\mathbf{A}} \mathbf{S} - \bar{\mathbf{A}}^T \bar{\mathbf{X}})}{\text{vec}(\mathbf{A}_k \mathbf{P}_S)^T \text{vec}(\mathbf{A}_k \mathbf{P}_S)} \quad (15)$$

are the adaptive steepest descent learning rates.

3. SIMULATION RESULTS

All the NMF algorithms discussed in this paper have been extensively tested for many difficult benchmarks for signals and images with various statistical distributions and also for real EEG data. We found the best performance can be obtained with the AIPG, FPALS and the algorithm (8)-(9) for $\beta = 1$.

Due to space limitation, we present here only one simulation Example: Six natural highly correlated images are mixed by randomly generated 3D tensor $\underline{\mathbf{A}} \in \mathbb{R}_+^{12 \times 6 \times 25}$. The observed mixed data are collected in 3D tensor $\underline{\mathbf{X}} \in \mathbb{R}_+^{12 \times 64^2 \times 25}$. The exemplary results are shown in Fig.2.

4. CONCLUSIONS AND DISCUSSION

In this paper we proposed generalized and flexible cost functions (controlled by a single parameter alpha or beta) that allows us to derive a family of robust and efficient NTF algorithms. The optimal choice of a free parameter of a specific cost function depends on a statistical distribution of data and additive noise, thus various criteria and algorithms (updating rules) should be applied for estimating the basis matrices \mathbf{A}_k and the source matrix \mathbf{S} , depending on *a priori* knowledge about the statistics of noise or errors. It is worth to mention that we can use three different strategies to estimate common factors (the source matrix \mathbf{S}). In the first approach, presented in this paper, we use two different cost functions: A global cost function (using unfolded column-wise matrices: $\bar{\mathbf{X}}$, $\bar{\mathbf{S}}$ for frontal slices of 3D tensors) to estimate the common factors \mathbf{S} , i.e., the source matrix \mathbf{S} ; and local cost functions to estimate the slices \mathbf{A}_k , ($k = 1, 2, \dots, K$). However, instead of using the unfolding matrices for the NTF2 model, in order to estimate \mathbf{S} , we can use, average matrices defined as $\bar{\mathbf{X}} = \sum_k \mathbf{X}_k \in \mathbb{R}^{I \times T}$ and $\bar{\mathbf{A}} = \sum_k \mathbf{A}_k \in \mathbb{R}^{I \times R}$. Furthermore, it is also possible to apply a different approach by using only set of local cost functions, e.g., $D_k(\mathbf{X}_k | \mathbf{A}_k \mathbf{S}) = 0.5 |\mathbf{X}_k - \mathbf{A}_k \mathbf{S}|_F^2$. In such a case, we estimate \mathbf{A}_k and \mathbf{S} cyclically by applying alternating minimization (similar to row-action projection of the Kaczmarz algorithm). We found that such approaches also work quite well for the NTF2 model. The advantage of the last approach is that the all updates learning rules are local (slice by slice) and algorithms are generally faster for large data, (especially, if $K \gg 1$).

Obviously, 3D NTF models can be transformed to a 2D non-negative matrix factorization (NMF) problem by unfold-

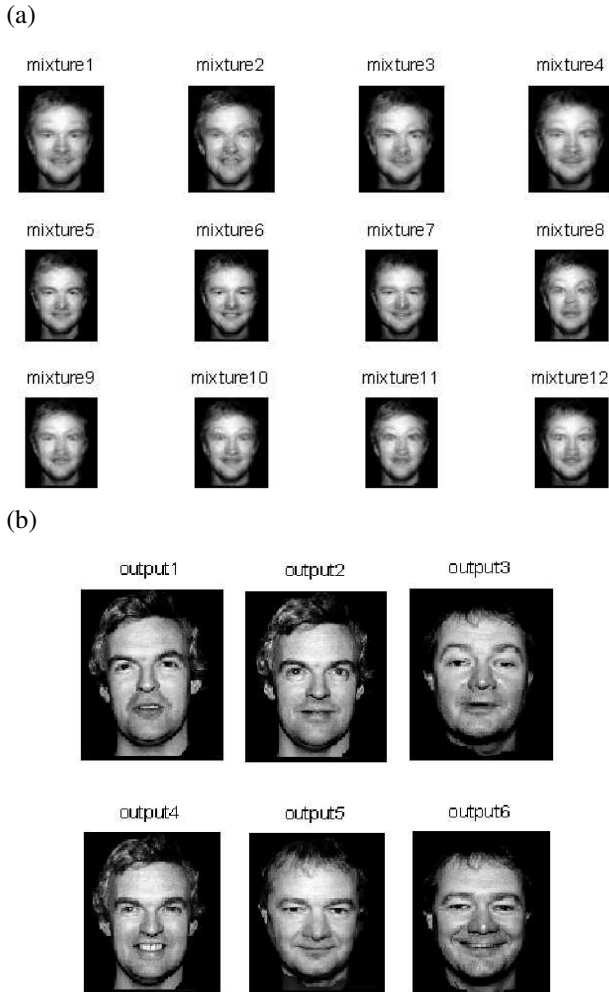


Fig. 2. Example 1: (a) Observed mixed images (only subset of the observed images are shown); (b) Estimated source images using the FPALS algorithm (SIR = 20.8, 22.7, 21.5, 23.1, 20.4, 24.2 [dB], respectively).

ing (matricizing) tensors. However, it should be noted that such a 2D model is not exactly equivalent to the NMF2 model, since in practice we often need to impose different additional constraints for each slice. In other words, the NTF2 model should not be considered as equal to a standard 2-way NMF of a single unfolded 2-D matrix. The profiles of the stacked (column-wise unfolded) \bar{A} are often not treated as single profiles and the constraints are usually applied independently to each A_k sub-matrix that form the stacked \bar{A} . Moreover, the NTF2 is considered as a dynamical process, where the data analysis is performed several times under different conditions (multi-start initializations, multilayer or recurrent implementation, Monte Carlo analysis, selection of additional natural constraints, etc.) to get full information about the available data and/or discover some inner structures in the data.

We have been motivated to develop proposed NTF algo-

gorithm by using them in three areas of data analysis (especially, EEG data) and signal/image processing: (i) multi-way blind source separation, (ii) model reductions and selection and (iii) sparse image coding. The proposed models can be further extended by imposing additional, natural constraints such as smoothness, continuity, closure, unimodality, local rank, selectivity, and/or by taking into account a prior knowledge about specific 3D, or more generally, multi-way data. Obviously, there are many challenging open issues remaining, such as global convergence, optimal choice of parameters and uniqueness of a solution when additional constraints are imposed.

5. REFERENCES

- [1] S. Amari, *Differential-Geometrical Methods in Statistics*, 1985 Springer Verlag.
- [2] M. Berry, M. Browne, A. Langville, P. Pauca, and R. Plemmons, "Algorithms for approximate nonnegative matrix factorization", *Computational Statistics and Data Analysis*, 2006 <http://www.wfu.edu/~plemmons/papers.htm>.
- [3] A. Cichocki, R. Zdunek, and S. Amari, "Csiszar's divergences for non-negative matrix factorization: Family of new algorithms" ICA-2006, Springer LNCS 3889, 32-39, 2006.
- [4] A. Cichocki, S. Amari, R. Zdunek, R. Kompass, G. Hori, and Z. He, "Extended SMART algorithms for non-negative matrix factorization", Springer LNAI 4029, 548-562 2006.
- [5] A. Cichocki and R. Zdunek, "NTFLAB for Signal and Image Processing", Tech. rep., Laboratory for Advanced Brain Signal Processing, BSI RIKEN, Saitama, Japan, <http://www.bsp.brain.riken.jp> (2006).
- [6] L. De Lathauwer, and P. Comon (Editors), Workshop on Tensor Decompositions and Applications, CIRM, Marseille, France, 2005 <http://www.etis.ensea.fr/wtda/>
- [7] I. Dhillon and S. Sra, "Generalized nonnegative matrix approximations with Bregman divergences". In: NIPS -Neural Information Proc. Systems, Vancouver Canada, 2005.
- [8] T. Hazan, S. Polak, and A. Shashua, "Sparse image coding using a 3D non-negative tensor factorization", in Proc. of ICCV, 50-57, 2005.
- [9] M. Heiler and C. Schnoerr, "Controlling sparseness in Non-Negative Tensor Factorization", ECCV 2006, <http://www.cvprp.uni-mannheim.de/heiler/>
- [10] P. Hoyer, "Non-negative matrix factorization with sparseness constraints" *Journal of Machine Learning Research* 5, 1457-1469, 2004
- [11] M. Kim, and S. Choi, "Monaural music source separation: Nonnegativity, sparseness, and shift-invariance", ICA-2006, Springer LNCS 3889, pp. 617-624, 2006.
- [12] D.D. Lee, and H.S. Seung, "Learning of the parts of objects by non-negative matrix factorization". *Nature*, 401, 788-791, 1999.
- [13] H. Li, T. Adali, and D.E. Wang, "Non-negative matrix factorization with orthogonality constraints for chemical agent detection in Raman spectra", In: IEEE Workshop on Machine Learning for Signal Processing, Mystic USA, 2005.
- [14] M. Merritt, and Y. Zhang, "Interior-point gradient method for large-scale totally nonnegative least squares problems" *Journal of Optimization Theory and Applications*, 126(1), 193-202, 2005.
- [15] M. Morup, L.K. Hansen, C.S. Herrmann, J. Parnas and S.M. Arnfred, "Parallel factor analysis as an exploratory tool for wavelet transformed event-related EEG", *NeuroImage*, 29(3):938-947, 2006.
- [16] F. Miwakeichi, E. Martinez-Montes, P.A. Valdes-Sosa, N. Nishiyama, H. Mizuhara, and Y. Yamaguchi, "Decomposing EEG data into spacetime frequency components using parallel factor analysis", *NeuroImage*, 22, 1035-1045, 2004.
- [17] A. Smilde, R. Bro, and P. Geladi, *Multi-way Analysis: Applications in the Chemical Sciences*, Wiley, NY, 2004.
- [18] M.A.O. Vasilescu and D. Terzopoulos, "Multilinear analysis of image ensembles: Tensorfaces" In Proc. European Conf. on Computer Vision (ECCV), pages 447-460, Copenhagen, Denmark, 2002.
- [19] S.A. Vorobyov, Y. Rong, N.D. Sidiropoulos, and A.B. Gershman, "Robust iterative fitting of multilinear models", *IEEE Trans. Signal Processing*, vol. 53, no. 8, pp. 2678-2689, 2005.
- [20] R. Kompass, "A generalized divergence measure for nonnegative matrix factorization", *Neural Computation*, 2006 (in print).