

# ROBUST ITERATIVE FITTING OF MULTILINEAR MODELS BASED ON LINEAR PROGRAMMING

Sergiy A. Vorobyov<sup>†</sup> Yue Rong<sup>†</sup> Nicholas D. Sidiropoulos<sup>‡</sup> Alex B. Gershman<sup>†</sup>

<sup>†</sup>Dept. of Communication Systems, University of Duisburg-Essen, Germany

<sup>‡</sup>Dept. of Electronic and Computer Engineering, Technical University of Crete, Chania, Greece

## ABSTRACT

PARALLEL FACTOR (PARAFAC) analysis is an extension of low-rank matrix decomposition to higher-way arrays. It decomposes a given array in a sum of multilinear terms. PARAFAC analysis generalizes and unifies common array processing models (like joint diagonalization and ESPRIT); it has found numerous applications from blind multiuser detection and multi-dimensional harmonic retrieval, to clustering and nuclear magnetic resonance. The prevailing fitting algorithm in all these applications is based on alternating least squares (ALS) optimization, which is matched to Gaussian noise. In many cases, however, measurement errors are far from being Gaussian. In this paper, we develop an iterative algorithm for least absolute error fitting of general multilinear models, based on efficient interior point methods for Linear Programming (LP). We also benchmark its performance in Laplacian, Cauchy, and Gaussian noise environments, versus the respective CRBs and the commonly used ALS algorithm.

## 1. INTRODUCTION

The PARAFAC model is a useful data analysis tool that has recently found applications in array signal processing and communications [1], [2]. When generalizing the concept of low-rank decomposition to higher-way arrays, PARAFAC is instrumental in the analysis of data arrays indexed by three or more independent variables, just like Singular Value Decomposition (SVD) is instrumental in ordinary matrix (two-way array) analysis. In most applications of PARAFAC analysis, the ALS regression procedure is used to fit the model parameters [1], [2]. Least Squares (LS) regression is optimal (in the maximum likelihood sense) when measurement errors are additive i.i.d. Gaussian. Gaussianity is an often-made assumption, due to the central limit theorem, but also for tractability considerations. However, in many applications the measurement errors are far from being Gaussian random variables [3].

A. B. Gershman is on leave from the Department of Electrical and Computer Engineering, McMaster University, Hamilton, Ontario, Canada.

The Least Absolute Error (LAE) criterion is often used as a robust alternative to LS. LAE regression is optimal (in the maximum likelihood sense) when measurement errors are additive i.i.d. Laplacian. The Laplacian distribution is more heavy-tailed than the Gaussian one; therefore, it is better suited to model impulsive noise and outliers. Another distribution commonly used for modelling impulsive noise is the Cauchy, and, more generally, the class of  $\alpha$ -stable distributions [3]. It is therefore of interest to develop PARAFAC regression procedures that optimize the LAE fitting criterion.

In this paper we develop such an iterative procedure that makes use of LP. The performance of the proposed algorithm is illustrated by means of simulations and compared to the pertinent Cramér-Rao bounds (CRBs) and Trilinear ALS (TALS) procedure [1].

## 2. PARALLEL FACTOR ANALYSIS

We introduce notation that will be useful in the sequel. Consider an  $I \times J \times K$  three-way array  $\underline{\mathbf{X}}$  with typical element  $x_{i,j,k}$  and the  $F$ -component trilinear decomposition

$$x_{i,j,k} = \sum_{f=1}^F a_{i,f} b_{j,f} c_{k,f} \quad (1)$$

for all  $i = 1, \dots, I$ ,  $j = 1, \dots, J$  and  $k = 1, \dots, K$ . Here  $a_{i,f}$  stands for the  $(i, f)$ -th element of  $I \times F$  matrix  $\mathbf{A}$ , and similarly  $b_{j,f}$  and  $c_{k,f}$  stand for  $(j, f)$ -th and  $(k, f)$ -th elements of  $J \times F$  and  $K \times F$  matrices  $\mathbf{B}$  and  $\mathbf{C}$ , respectively. Matrices  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{C}$  are in general complex-valued. Equation (1) expresses  $x_{i,j,k}$  as a sum of  $F$  rank-one triple products; it is known as trilinear decomposition, or *PARAFAC analysis* model of  $x_{i,j,k}$ .

Let  $\mathbf{A}_i = \mathcal{D}_i(\mathbf{A})$  denote the operator which takes the  $i$ -th row of matrix  $\mathbf{A}$  and produces a diagonal matrix by placing this row on the main diagonal. Then by “slicing” the three-dimensional array  $\underline{\mathbf{X}}$  in a series of “slabs” (two-dimensional arrays), we obtain

$$\mathbf{X}_i = \mathbf{B} \mathbf{A}_i \mathbf{C}^T, \quad i = 1, \dots, I \quad (2)$$

Here such a slicing is made perpendicular to the  $i$ th dimension, i.e.,  $\mathbf{X}_i := [x_{i,\cdot}]$  is the  $J \times K$  two-dimensional slice of  $\underline{\mathbf{X}}$  corresponding to the given index  $i$ . Two other types of slicing of  $\underline{\mathbf{X}}$  are useful in understanding the algorithm that will be developed in the next section. They are given by

$$\mathbf{Y}_j = \mathbf{C}\mathbf{B}_j\mathbf{A}^T, \quad j = 1, \dots, J \quad (3)$$

$$\mathbf{Z}_k = \mathbf{A}\mathbf{C}_k\mathbf{B}^T, \quad k = 1, \dots, K \quad (4)$$

where  $\mathbf{B}_j = \mathcal{D}_j(\mathbf{B})$ ,  $\mathbf{C}_k = \mathcal{D}_k(\mathbf{C})$ , while the  $K \times I$  matrix  $\mathbf{Y}_j$  and  $I \times J$  matrix  $\mathbf{Z}_k$  are defined as  $\mathbf{Y}_j := [x_{\cdot,j,\cdot}]$  and  $\mathbf{Z}_k := [x_{\cdot,\cdot,k}]$ , respectively.

### 3. TRILINEAR ALTERNATING LAE REGRESSION BASED ON LINEAR PROGRAMMING (TALAE-LP)

In practice, the three-way array will contain measurement noise, i.e.  $\widetilde{\mathbf{X}} = \mathbf{X} + \mathbf{V}$  where the  $(i, j, k)$ th element of  $\widetilde{\mathbf{X}}$  can be written as

$$\tilde{x}_{i,j,k} = x_{i,j,k} + v_{i,j,k} \quad (5)$$

and  $v_{i,j,k}$  denotes the additive complex i.i.d. zero-mean measurement noise with statistically independent real and imaginary parts.

The PARAFAC fitting problem is then formulated as follows. We are given the noisy data  $\widetilde{\mathbf{X}}$  and wish to estimate  $\mathbf{A}$ ,  $\mathbf{B}$ , and  $\mathbf{C}$ . Let us introduce the tall matrix

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_I \end{bmatrix}_{JI \times K} = \begin{bmatrix} \mathbf{B}\mathbf{A}_1 \\ \vdots \\ \mathbf{B}\mathbf{A}_I \end{bmatrix} \mathbf{C}^T = (\mathbf{A} \odot \mathbf{B})\mathbf{C}^T \quad (6)$$

where  $\odot$  stands for the Khatri-Rao matrix product. Similarly, we introduce the matrix of noisy data

$$\widetilde{\mathbf{X}} = \begin{bmatrix} \widetilde{\mathbf{X}}_1 \\ \vdots \\ \widetilde{\mathbf{X}}_I \end{bmatrix} = \begin{bmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_I \end{bmatrix} + \begin{bmatrix} \mathbf{V}_1 \\ \vdots \\ \mathbf{V}_I \end{bmatrix} \quad (7)$$

Then the conditional Maximum Likelihood (ML) estimation problem for the matrix  $\mathbf{C}$  given matrices  $\mathbf{A}$  and  $\mathbf{B}$  and assuming i.i.d. Gaussian measurement noise is the LS fitting problem  $\min_{\mathbf{C}} \|\widetilde{\mathbf{X}} - (\mathbf{A} \odot \mathbf{B})\mathbf{C}^T\|_F^2$  where  $\|\cdot\|_F^2$  denotes the Frobenius matrix norm.

If the measurement noise is i.i.d. Laplacian (with i.i.d. Laplacian-distributed real and imaginary parts in the complex case), then ML estimation is equivalent to LAE regression. Some manipulations are necessary in order to express the absolute error criterion in the form of a convenient vector  $\ell_1$  norm. Towards this end, introduce the operator  $\mathcal{F}(\cdot)$

$$\mathbf{s} = \mathcal{F}(\mathbf{S}) = \begin{bmatrix} \check{\mathbf{S}}_{\cdot,1} \\ \vdots \\ \check{\mathbf{S}}_{\cdot,L} \end{bmatrix}, \quad \check{\mathbf{S}}_{\cdot,l} = \begin{bmatrix} \text{Re}\{\mathbf{S}_{\cdot,l}\} \\ \text{Im}\{\mathbf{S}_{\cdot,l}\} \end{bmatrix} \quad (8)$$

where  $\mathbf{S}$  is a complex-valued  $M \times L$  matrix, and  $\mathbf{S}_{\cdot,l}$  denotes its  $l$ th column. The following property holds:

$$\mathcal{F}\{\mathbf{D}\mathbf{F}\} = (\mathbf{I} \otimes \mathcal{G}\{\mathbf{D}\})\mathcal{F}\{\mathbf{F}\} \quad (9)$$

where  $\mathbf{I}$  is the identity matrix of a commensurate dimension,  $\mathbf{D}$  and  $\mathbf{F}$  are any complex-valued matrices of commensurate dimensions,  $\otimes$  denotes the Kronecker matrix product, and  $\mathcal{G}\{\mathbf{D}\}$  stands for the operator

$$\mathcal{G}\{\mathbf{D}\} = \begin{bmatrix} \text{Re}\{\mathbf{D}\} & -\text{Im}\{\mathbf{D}\} \\ \text{Im}\{\mathbf{D}\} & \text{Re}\{\mathbf{D}\} \end{bmatrix} \quad (10)$$

Using property (9), we find that the absolute error model fitting criterion can be written as

$$\|\tilde{\mathbf{x}} - (\mathbf{I}_K \otimes \mathcal{G}\{\mathbf{A} \odot \mathbf{B}\})\mathbf{c}\|_1 \quad (11)$$

i.e., through the  $\ell_1$  norm of a real-valued vector. Here,  $\tilde{\mathbf{x}} = \mathcal{F}(\widetilde{\mathbf{X}})$ ,  $\mathbf{c} = \mathcal{F}(\mathbf{C})$ , and the dimension of the identity matrix is clarified by means of the subscript  $K$ .

Using the other two ways of slicing the array  $\underline{\mathbf{X}}$ , we introduce the matrices  $\mathbf{Y} = (\mathbf{B} \odot \mathbf{C})\mathbf{A}^T$  and  $\mathbf{Z} = (\mathbf{C} \odot \mathbf{A})\mathbf{B}^T$ . Furthermore, we introduce

$$\widetilde{\mathbf{Y}} = \begin{bmatrix} \widetilde{\mathbf{Y}}_1 \\ \vdots \\ \widetilde{\mathbf{Y}}_J \end{bmatrix}_{KJ \times I}, \quad \widetilde{\mathbf{Z}} = \begin{bmatrix} \widetilde{\mathbf{Z}}_1 \\ \vdots \\ \widetilde{\mathbf{Z}}_K \end{bmatrix}_{IK \times J} \quad (12)$$

where  $\widetilde{\mathbf{Y}}_j$ ,  $j = 1, \dots, J$ , and  $\widetilde{\mathbf{Z}}_k$ ,  $k = 1, \dots, K$  are the noisy slabs of  $\widetilde{\mathbf{X}}$  along corresponding dimensions.

Now we have all notations necessary to explain the new fitting algorithm.

The idea behind this algorithm is similar to that of TALS regression for Gaussian noise [2] and is as follows: each time, update a subset of parameters using the LAE criterion, conditioned on previously obtained estimates of the remaining parameters; proceed to update another subset of parameters; repeat until convergence.

In more detail, we first initialize matrices  $\mathbf{A}$  and  $\mathbf{B}$  randomly or by single-invariance ESPRIT when applicable [1], [2]. Then, given the matrix  $\widetilde{\mathbf{X}}$ , and these initial estimates of  $\mathbf{A}$  and  $\mathbf{B}$  (which we denote hereafter as  $\widehat{\mathbf{A}}$  and  $\widehat{\mathbf{B}}$ ), our purpose is to find the estimate of the matrix  $\mathbf{C}$  which minimizes the norm (11). Specifically, we have to find the estimate of  $\mathbf{C}$  by solving the following optimization problem:

$$\begin{aligned} \hat{\mathbf{c}} &= \arg \min_{\mathbf{c}} \|\tilde{\mathbf{x}} - (\mathbf{I}_K \otimes \mathcal{G}\{\widehat{\mathbf{A}} \odot \widehat{\mathbf{B}}\})\mathbf{c}\|_1 \\ \widehat{\mathbf{C}} &= (\mathcal{F}^{-1}\{\hat{\mathbf{c}}\})^T \end{aligned} \quad (13)$$

for given  $\tilde{\mathbf{x}}$ ,  $\widehat{\mathbf{A}}$  and  $\widehat{\mathbf{B}}$ . In (13),  $\mathcal{F}^{-1}\{\cdot\}$  denotes the inverse operator to  $\mathcal{F}\{\cdot\}$  of (8). Introducing the vector  $\mathbf{e} = [1, 1, \dots, 1]^T$  and the vector of slack variables  $\mathbf{q}_1$  (both of

commensurate dimensions), we can equivalently write the problem (13) as

$$\begin{aligned} \min_{c, q_1} e^T q_1 \quad \text{subject to} \quad & \tilde{x} - (\mathbf{I}_K \otimes \mathcal{G}\{\hat{\mathbf{A}} \odot \hat{\mathbf{B}}\})c \preceq q_1 \\ & \tilde{x} - (\mathbf{I}_K \otimes \mathcal{G}\{\hat{\mathbf{A}} \odot \hat{\mathbf{B}}\})c \succeq -q_1 \end{aligned}$$

where  $\succeq$  denotes the usual *pointwise* ordering. This optimization problem is an LP problem that can be very efficiently solved using interior-point methods [4].

Using the second way of slicing the three-dimensional array (i.e., working with the data  $\tilde{\mathbf{y}} = \mathcal{F}(\tilde{\mathbf{Y}})$ ) and exploiting the property (9), we obtain that the estimate of  $\mathbf{A}$  can be found by solving the following optimization problem:

$$\begin{aligned} \hat{\mathbf{a}} &= \arg \min_{\mathbf{a}} \|\tilde{\mathbf{y}} - (\mathbf{I}_I \otimes \mathcal{G}\{\hat{\mathbf{B}} \odot \hat{\mathbf{C}}\})\mathbf{a}\|_1 \\ \hat{\mathbf{A}} &= (\mathcal{F}^{-1}\{\hat{\mathbf{a}}\})^T \end{aligned} \quad (14)$$

with given  $\tilde{\mathbf{y}}$  and previously estimated  $\hat{\mathbf{B}}$  and  $\hat{\mathbf{C}}$ . This problem can be rewritten as the following LP problem:

$$\begin{aligned} \min_{\mathbf{a}, q_2} e^T q_2 \quad \text{subject to} \quad & \tilde{\mathbf{y}} - (\mathbf{I}_I \otimes \mathcal{G}\{\hat{\mathbf{B}} \odot \hat{\mathbf{C}}\})\mathbf{a} \preceq q_2 \\ & \tilde{\mathbf{y}} - (\mathbf{I}_I \otimes \mathcal{G}\{\hat{\mathbf{B}} \odot \hat{\mathbf{C}}\})\mathbf{a} \succeq -q_2 \end{aligned}$$

where  $q_2$  is the vector of slack variables of commensurate dimension.

Finally, using the third way of slicing the three-dimensional array and applying the property (9), we can find the estimate of  $\mathbf{B}$  by solving the following optimization problem:

$$\begin{aligned} \hat{\mathbf{b}} &= \arg \min_{\mathbf{b}} \|\tilde{\mathbf{z}} - (\mathbf{I}_J \otimes \mathcal{G}\{\hat{\mathbf{C}} \odot \hat{\mathbf{A}}\})\mathbf{b}\|_1 \\ \hat{\mathbf{B}} &= (\mathcal{F}^{-1}\{\hat{\mathbf{b}}\})^T \end{aligned} \quad (15)$$

with given  $\tilde{\mathbf{z}}$  and previously estimated  $\hat{\mathbf{A}}$  and  $\hat{\mathbf{C}}$ . This problem is equivalent to the following LP problem:

$$\begin{aligned} \min_{\mathbf{b}, q_3} e^T q_3 \quad \text{subject to} \quad & \tilde{\mathbf{z}} - (\mathbf{I}_J \otimes \mathcal{G}\{\hat{\mathbf{C}} \odot \hat{\mathbf{A}}\})\mathbf{b} \preceq q_3 \\ & \tilde{\mathbf{z}} - (\mathbf{I}_J \otimes \mathcal{G}\{\hat{\mathbf{C}} \odot \hat{\mathbf{A}}\})\mathbf{b} \succeq -q_3 \end{aligned}$$

where  $q_3$  is the vector of slack variables of commensurate dimension.

Fitting proceeds by updating one matrix at a time, conditioned on interim estimates of the other two, in a round-robin fashion. Note that the conditional update of any given matrix may either improve or maintain but cannot worsen the current fit. Monotone convergence of the fit (but not necessarily to the global minimum) follows directly from this observation. The per-iteration complexity of TALAE-LP is equal to the cost of solving LP problems [4]. This is of the same order of complexity as computing a matrix pseudo-inverse in the TALS method [1] and can be estimated as  $\mathcal{O}(F^3 + FIJK)$ . Overall complexity depends on the number of iterations, which varies depending on problem-specific parameters and the given batch of data.

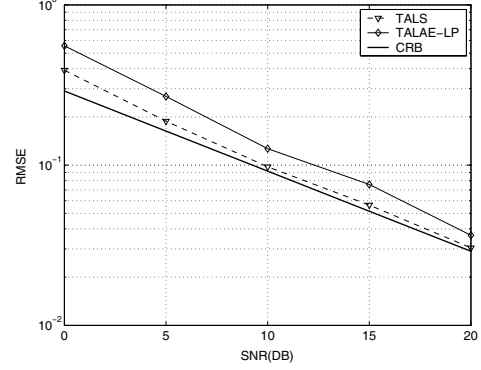


Figure 1: RMSEs versus SNR. Gaussian channel noise.

#### 4. SIMULATIONS

Let us compare the performance of the proposed TALAE-LP algorithm with that of the conventional TALS method, and against the pertinent CRB. The example of blind PARAFAC multiuser detection for a Direct-Sequence Code Division Multiple Access (DS-CDMA) communication system [1] is simulated. For the DS-CDMA application, the elements of the matrices  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{C}$  have the following meanings:  $a_{i,f}$  is fading/gain between user  $f$  and antenna element  $i$ ;  $b_{j,f}$  is  $j$ -th chip of the spreading code of user  $f$ ;  $c_{k,f}$  is  $k$ -th symbol transmitted by user  $f$ . Correspondingly, matrix  $\mathbf{A}$  is the channel matrix,  $\mathbf{B}$  is the spreading code matrix, and  $\mathbf{C}$  is the user symbol matrix, all unknown to the receiver. Here,  $F$  is the number of users,  $I$  is the total number of antenna elements in the array,  $J$  is the number of Intersymbol Interference (ISI)-free chips per symbol, and  $K$  is the length of the transmitted sequence of symbols.

The data  $\mathbf{X}$  are contaminated by channel noise. Three models of the channel noise are used. One is Gaussian noise, while the other two are Laplacian and Cauchy noise.

For LS fitting, we use the COMFAC algorithm [1] which is a fast implementation of TALS. The MOSEK convex optimization MATLAB toolbox [5] is used to solve LP problems associated with our TALAE-LP algorithm. Scale and permutation ambiguities are inherent to this blind separation problem [1]; the scale ambiguity manifests itself as a complex constant that multiplies each individual row of  $\mathbf{C}$ . For constant-modulus transmissions, this ambiguity can be removed via Automatic Gain Control (AGC) and differential encoding/decoding. We assume differentially-encoded user signals. For the purpose of performance evaluation only, the permutation ambiguity is resolved using a greedy LS ( $\mathbf{C}, \hat{\mathbf{C}}$ ) column-matching algorithm.

We present Monte Carlo simulations that are designed to assess the Root Mean Square Error (RMSE) performance of the aforementioned algorithms. The parameters used in the simulations are as follows:  $N$  = number of Monte Carlo tri-

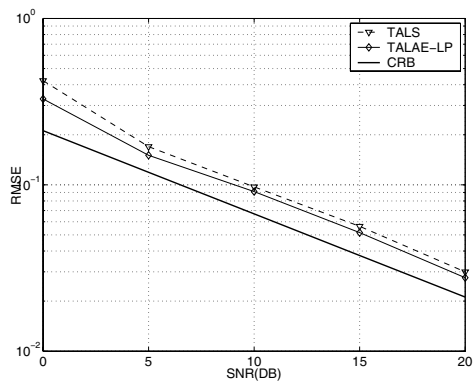


Figure 2: RMSEs versus SNR. Laplacian channel noise.

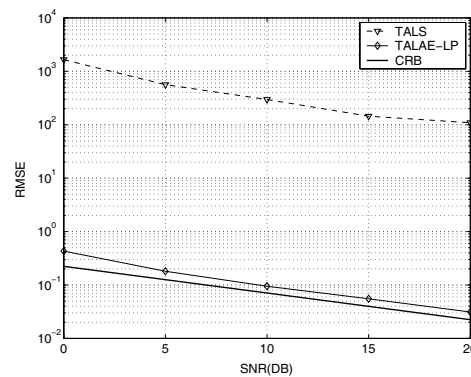


Figure 3: RMSEs versus SNR. Cauchy channel noise.

als = 100;  $F = 2$ ;  $I = 8$ ;  $J = 8$ ;  $K = 20$ ; and  $\alpha = 1$ . Here,  $\alpha$  is the characteristic exponent which determines the heaviness of the tail of the symmetric  $\alpha$ -stable distribution used in our third example ( $\alpha = 1$  yields the Cauchy distribution). The associated symmetric  $\alpha$ -stable characteristic function is given by  $\phi(\omega) = \exp\{-\gamma|\omega|^\alpha\}$ , where  $\gamma$  is a positive constant related to the scale of the distribution. The *geometric* Signal-to-Noise Ratio (SNR) in this case is defined as [6]  $\text{SNR}_{\text{Cauchy}} = \frac{A^2}{4C_g\gamma^2}$ , where  $C_g = e^{C_e} \approx 1.78$  and  $A$  is the magnitude of the noise-free signal. For the Gaussian case, the geometric SNR is equivalent to the standard SNR, and in the Laplacian case we use the standard SNR, since Laplacian distribution does not belong to the class of symmetric  $\alpha$ -stable distributions.

Throughout the simulations, we assume that the noise power is normalized to be equal to 1. User signals are re-drawn from an i.i.d. Bernoulli distribution and differentially encoded for each Monte Carlo trial. BPSK modulation is used for all user signals. The gains of the channel matrix  $\mathbf{A}$  and the elements of the spreading code matrix  $\mathbf{B}$  are generated as i.i.d. Gaussian unit variance random variables and are fixed in each Monte Carlo trial, while re-drawn from one trial to another.

Figures 1, 2 and 3 display the performance of the aforementioned algorithms in terms of RMSE versus the SNR for the case of Gaussian, Laplacian and Cauchy noise, respectively, and compare the performance with the corresponding CRBs. Figures 1 and 2 demonstrate that in the case of Gaussian noise, the TALS method performs slightly better than the proposed robust TALAE-LP algorithm, while in the case of Laplacian noise, the TALAE-LP algorithm has slightly better performance as compared to the TALS method. In the case of Cauchy noise (Fig. 3), the TALS method breaks down, while the TALAE-LP algorithm is not affected and is close to the CRB. The degradation in performance relative to TALS in the Gaussian case can be considered as a moderate price for greatly improved robustness against heavy-tailed Cauchy noise.

## 5. CONCLUSIONS

An iterative algorithm for robust fitting of trilinear PARAFAC models has been proposed. The algorithm relies on alternating optimization using LP. The proposed algorithm outperforms the popular alternative LS PARAFAC fitting procedure under heavy-tailed noise. Even though our algorithm is matched to the Laplacian distribution, it performs very well under Cauchy noise. Furthermore, its performance degrades only moderately under Gaussian noise.

## 6. REFERENCES

- [1] N. D. Sidiropoulos, G. B. Giannakis and R. Bro, "Blind PARAFAC receivers for DS-CDMA systems," *IEEE Trans. Signal Processing*, vol. 48, pp. 810-823, Mar. 2000.
- [2] N. D. Sidiropoulos, R. Bro and G. B. Giannakis, "Parallel factor analysis in sensor array processing," *IEEE Trans. Signal Processing*, vol. 48, pp. 2377-2388, Aug. 2000.
- [3] P. Tsakalides and C.L. Nikias, "Maximum likelihood localization of sources in noise modeled as a stable process," *IEEE Trans. Signal Processing*, vol. 43, pp. 2700-2713, Nov. 1995.
- [4] Y. Nesterov and A. Nemirovski, *Interior Point Polynomial Algorithms in Convex Programming*. Philadelphia, PA: SIAM, 1994.
- [5] *The MOSEK optimization tools version 2.0 (Build 19). User's manual and reference*, EKA Consulting ApS, Denmark, 2001. <http://www.mosek.com>
- [6] J. G. Gonzalez, "Robust techniques for wireless communications in non-Gaussian environments," Ph.D. dissertation, University of Delaware, 1997.