



HAL
open science

Blind audio source separation using sparsity based criterion for convolutive mixture case

Abdeldjalil Aissa El Bey, Karim Abed-Meraim, Yves Grenier

► **To cite this version:**

Abdeldjalil Aissa El Bey, Karim Abed-Meraim, Yves Grenier. Blind audio source separation using sparsity based criterion for convolutive mixture case. Independent Component Analysis and Signal Separation (ICA), Sep 2007, Londres, United Kingdom. pp.317-324. hal-01772799

HAL Id: hal-01772799

<https://hal.science/hal-01772799v1>

Submitted on 20 Apr 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Blind audio source separation using sparsity based criterion for convolutive mixture case

A. Aïssa-El-Bey, K. Abed-Meraim and Y. Grenier

ENST-Paris, TSI Department, 46 rue Barrault 75634, Paris Cedex 13, France
{elbey, abed, grenier}@tsi.enst.fr

Abstract. In this paper, we are interested in the separation of audio sources from their instantaneous or convolutive mixtures. We propose a new separation method that exploits the sparsity of the audio signals via an ℓ_p -norm based contrast function. A simple and efficient natural gradient technique is used for the optimization of the contrast function in an instantaneous mixture case. We extend this method to the convolutive mixture case, by exploiting the property of the Fourier transform. The resulting algorithm is shown to outperform existing techniques in terms of separation quality and computational cost.

1 Introduction

Blind Source Separation (BSS) is an approach to estimate and recover independent source signals using only the information within the mixtures observed at each channel. Many algorithms have been proposed to solve the standard blind source separation problem in which the mixtures are assumed to be instantaneous. A fundamental and necessary assumption of BSS is that the sources are statistically independent and thus are often separated using higher-order statistical information [1]. If extra information about the sources is available at hand, such as temporal coherency [2], source nonstationarity [3], or source cyclostationarity [4], then one can remain in the second-order statistical scenario, to achieve the BSS.

In the case of non-stationary signals (including audio signals), certain solutions using time-frequency analysis of the observations exist [5]. Other solutions use the statistical independence of the sources assuming a local stationarity to solve the BSS problem [6]. This is a strong assumption that is not always verified [7]. To avoid this problem, we propose a new approach that handles the general linear instantaneous model (possibly noisy) by using the *sparsity* assumption of the sources in the time domain. Then, we extend this algorithm to the convolutive mixture case, by transforming the convolutive problem into instantaneous problem in the frequency domain, and separating the instantaneous mixtures in every frequency bin. The use of sparsity to handle this model, has arisen in several papers in the area of source separation [8, 9]. We first present a sparsity contrast function for BSS. Then, in order to achieve BSS, we optimize the considered contrast function using an iterative algorithm based on the relative gradient technique.

In the following section, we discuss the data model that formulates our problem. Next, we detail the different steps of the proposed algorithm. In Section 4, some simulations are undertaken to validate our algorithm and to compare its performance to other existing BSS techniques.

2 Instantaneous mixture case

2.1 Data model

Assume that N audio signals impinge on an array of $M \geq N$ sensors. The measured array output is a weighted superposition of the signals, corrupted by additive noise, i.e.

$$\mathbf{x}(t) = \mathbf{A}\mathbf{s}(t) + \mathbf{w}(t) \quad t = 0, \dots, T-1 \quad (1)$$

where $\mathbf{s}(t) = [s_1(t), \dots, s_N(t)]^T$ is the $N \times 1$ sparse source vector, $\mathbf{w}(t) = [w_1(t), \dots, w_M(t)]^T$ is the $M \times 1$ complex noise vector, \mathbf{A} is the $M \times N$ full column rank mixing matrix (i.e., $M \geq N$), and the superscript T denotes the transpose operator. The purpose of blind source separation is to find a separating matrix, i.e. a $N \times M$ matrix \mathbf{B} such that $\hat{\mathbf{s}}(t) = \mathbf{B}\mathbf{x}(t)$ is an estimate of the source signals.

Before proceeding, note that complete blind identification of separating matrix \mathbf{B} (or equivalently, the mixing matrix \mathbf{A}) is impossible in this context, because the exchange of a fixed scalar between the source signal and the corresponding column of \mathbf{A} leaves the observations unaffected. Also note that the *numbering* of the signals is immaterial. It follows that the best that can be done is to determine \mathbf{B} up to a permutation and scalar shifts of its columns, i.e., \mathbf{B} is a separating matrix iff :

$$\mathbf{B}\mathbf{x}(t) = \mathbf{P}\mathbf{A}\mathbf{s}(t) \quad (2)$$

where \mathbf{P} is a permutation matrix and \mathbf{A} a non-singular diagonal matrix.

2.2 Sparsity-based BSS algorithm

Before starting, we propose to use 'an optional' whitening step which set the mixtures to the same energy level and reduces the number of parameters to be estimated. More precisely, the whitening step is applied to the signal mixtures before using our separation algorithm. The whitening is achieved by applying a $N \times M$ matrix \mathbf{W} to the signal mixtures in such a way $\text{Cov}(\mathbf{W}\mathbf{x}) = \mathbf{I}$ in the noiseless case, where $\text{Cov}(\cdot)$ stands for the covariance operator. As shown in [2], \mathbf{W} can be computed as the inverse square root of the noiseless covariance matrix of the signal mixtures (see [2] for more details). In the following, we apply our separation algorithm on the whitened data :

$$\mathbf{x}_w(t) = \mathbf{W}\mathbf{x}(t).$$

We propose an iterative algorithm for the separation of sparse audio signals, namely the ISBS for Iterative Sparse Blind Separation. It is well known that

audio signals are characterized by their sparsity property in the time domain [8, 9] which is measured by their ℓ_p norm where $0 \leq p < 2$. More specifically, one can define the following sparsity based contrast function

$$G_p(\mathbf{s}) = \sum_{i=1}^N [\mathcal{J}_p(s_i)]^{\frac{1}{p}} , \quad (3)$$

where

$$\mathcal{J}_p(s_i) = \frac{1}{T} \sum_{t=0}^{T-1} |s_i(t)|^p . \quad (4)$$

The algorithm finds a separating matrix \mathbf{B} such as,

$$\mathbf{B} = \arg \min_{\mathbf{B}} \{G_p(\mathbf{B})\} , \quad (5)$$

where

$$\mathcal{G}_p(\mathbf{B}) \triangleq G_p(\mathbf{z}) , \quad (6)$$

and $\mathbf{z}(t) \triangleq \mathbf{B}\mathbf{x}_w(t)$ represents the estimated sources. The approach we choose to solve (5) is inspired from [10]. It is a block technique based on the processing of T received samples and consists in searching iteratively the minimum of (5) in the form :

$$\mathbf{B}^{(k+1)} = (\mathbf{I} + \boldsymbol{\epsilon}^{(k)})\mathbf{B}^{(k)} \quad (7)$$

$$\mathbf{z}^{(k+1)}(t) = (\mathbf{I} + \boldsymbol{\epsilon}^{(k)})\mathbf{z}^{(k)}(t) \quad (8)$$

where \mathbf{I} denotes the identity matrix. At iteration k , a matrix $\boldsymbol{\epsilon}^{(k)}$ is determined from a local linearization of $G_p(\mathbf{B}^{(k+1)}\mathbf{x}_w)$. It is an approximate Newton technique with the benefit that $\boldsymbol{\epsilon}^{(k)}$ can be very simply computed (no Hessian inversion) under the additional assumption that $\mathbf{B}^{(k)}$ is close to a separating matrix. This procedure is illustrated in the following :

At the $(k+1)^{th}$ iteration, the proposed criterion (4) can be developed as follows:

$$\begin{aligned} \mathcal{J}_p(z_i^{(k+1)}) &= \frac{1}{T} \sum_{t=0}^{T-1} \left| z_i^{(k)}(t) + \sum_{j=1}^N \epsilon_{ij}^{(k)} z_j^{(k)}(t) \right|^p \\ &= \frac{1}{T} \sum_{t=0}^{T-1} |z_i^{(k)}(t)|^p \left| 1 + \sum_{j=1}^N \epsilon_{ij}^{(k)} \frac{z_j^{(k)}(t)}{z_i^{(k)}(t)} \right|^p . \end{aligned}$$

Under the assumption that $\mathbf{B}^{(k)}$ is close to a separating matrix, we have

$$|\epsilon_{ij}^{(k)}| \ll 1$$

and thus, a first order approximation of $\mathcal{J}_p(z_i^{(k+1)})$ is given by :

$$\begin{aligned} \mathcal{J}_p(z_i^{(k+1)}) &\approx \frac{1}{T} \sum_{t=0}^{T-1} |z_i^{(k)}(t)|^p + p \sum_{j=1}^N \Re(\epsilon_{ij}^{(k)}) \Re e \left(|z_i^{(k)}(t)|^{p-1} e^{-j\phi_i^{(k)}(t)} z_j^{(k)}(t) \right) \\ &\quad - \Im(\epsilon_{ij}^{(k)}) \Im m \left(|z_i^{(k)}(t)|^{p-1} e^{-j\phi_i^{(k)}(t)} z_j^{(k)}(t) \right) \end{aligned} \quad (9)$$

where $\Re e(x)$ and $\Im m(x)$ denote the real and imaginary parts of x and $\phi_i^{(k)}(t)$ is the argument of the complex number $z_i^{(k)}(t)$.

Using equation (9), equation (3) can be rewritten in more compact form as :

$$\mathcal{G}_p(\mathbf{B}^{(k+1)}) = \mathcal{G}_p(\mathbf{B}^{(k)}) + \Re e \left\{ \text{Tr} \left(\overline{\boldsymbol{\epsilon}}^{(k)} \mathbf{R}^{(k)H} \mathbf{D}^{(k)H} \right) \right\} \quad (10)$$

where $\overline{(\cdot)}$ denotes the conjugate of (\cdot) , $\text{Tr}(\cdot)$ is the matrix trace operator and the ij^{th} entry of matrix $\mathbf{R}^{(k)}$ is given by :

$$\mathcal{R}_{ij}^{(k)} = \frac{1}{T} \sum_{t=0}^{T-1} |z_i^{(k)}(t)|^{p-1} e^{-j\phi_i^{(k)}(t)} z_j^{(k)}(t) \quad (11)$$

$$\mathbf{D}^{(k)} = \left[\text{diag} \left(\mathcal{R}_{11}^{(k)}, \dots, \mathcal{R}_{NN}^{(k)} \right) \right]^{\frac{1}{p}-1}. \quad (12)$$

Using a gradient technique, $\boldsymbol{\epsilon}^{(k)}$ can be chosen as :

$$\boldsymbol{\epsilon}^{(k)} = -\mu \mathbf{D}^{(k)} \overline{\mathbf{R}}^{(k)} \quad (13)$$

where $\mu > 0$ is the gradient step. Replacing (13) into (10) leads to,

$$\mathcal{G}_p(\mathbf{B}^{(k+1)}) = \mathcal{G}_p(\mathbf{B}^{(k)}) - \mu \|\mathbf{D}^{(k)} \mathbf{R}^{(k)}\|^2. \quad (14)$$

So μ controls the decrement of the criterion. Now, to avoid the algorithm's convergence to the trivial solution $\mathbf{B} = \mathbf{0}$, one normalizes the outputs of the separating matrix to unit-power, i.e. $\rho_{z_i}^{(k+1)} \triangleq \frac{1}{T} \sum_{t=0}^{T-1} |z_i^{(k+1)}(t)|^2 = 1, \forall i$. Using first order approximation, this normalization leads to :

$$\epsilon_{ii}^{(k)} = \frac{1 - \rho_{z_i}^{(k)}}{2\rho_{z_i}^{(k)}}. \quad (15)$$

After convergence of the algorithm, the separation matrix $\mathbf{B} = \mathbf{B}^{(\mathcal{K})}$ is applied to the whitened signal mixtures \mathbf{x}_w to obtain an estimation of the original source signals. \mathcal{K} denotes here the number of iterations that can be either chosen a priori or given by a stopping criterion of the form $\|\mathbf{B}^{(k+1)} - \mathbf{B}^{(k)}\| < \delta$ where δ is a small threshold value.

3 Convulsive mixture case

Unfortunately, instantaneous mixing is very rarely encountered in real-world situations, where multipath propagation with large channel delay spread occurs, in which case convulsive mixtures are considered. In this case, the signal can be modeled by the following equation :

$$\mathbf{x}(t) = \sum_{l=0}^L \mathbf{H}(l)\mathbf{s}(t-l) + \mathbf{w}(t) \quad (16)$$

where $\mathbf{H}(l)$ are $M \times N$ matrices for $l \in [0, L]$ representing the impulse response coefficients of the channel and the polynomial matrix $\mathbf{H}(z) = \sum_{l=0}^L \mathbf{H}(l)z^{-l}$ is assumed to be irreducible (i.e. $\mathbf{H}(z)$ is of full column rank for all z). If we apply a short time Fourier transform (STFT) to the observed data $\mathbf{x}(t)$, the model in (16) (in the noiseless case) becomes approximately

$$\mathcal{S}_{\mathbf{x}}(t, f) \approx \mathbf{H}(f)\mathcal{S}_{\mathbf{s}}(t, f) \quad (17)$$

where $\mathcal{S}_{\mathbf{x}}(t, f)$ is the mixture STFT vector, $\mathcal{S}_{\mathbf{s}}(t, f)$ is the source STFT vector and $\mathbf{H}(f)$ is the channel Fourier Transform matrix. It shows that, for each frequency bin, the convulsive mixtures reduce to simple instantaneous mixtures. Therefore we can apply our ISBS algorithm for each frequency and separate the signals. As a result, in each frequency bin, we obtain the STFT source estimate

$$\mathcal{S}_{\hat{\mathbf{s}}}(t, f) = \mathbf{B}(f)\mathcal{S}_{\mathbf{x}}(t, f) . \quad (18)$$

It seems natural to reconstruct the separated signals by aligning these $\mathcal{S}_{\hat{\mathbf{s}}}(t, f)$ obtained for each frequency bin and applying the inverse short time Fourier transform. For that we need first to solve the permutation and scaling ambiguities as shown next.

3.1 Removing the scaling and permutation ambiguities

In this stage, the output of the separation filter is processed with the permutation matrix $\mathbf{\Pi}(f)$ and the scaling matrix $\mathbf{C}(f)$.

$$\mathbf{G}(f) = \mathbf{\Pi}(f)\mathbf{C}(f)\mathbf{B}(f) . \quad (19)$$

The scaling matrix $\mathbf{C}(f)$ is a $N \times N$ diagonal matrix found as in [11] by $\mathbf{C}(f) = \text{diag}[\mathbf{B}(f)^{\#}]$. For the permutation matrix $\mathbf{\Pi}(f)$, we exploit the continuity property of the acoustic filter in the frequency domain [12]. To align the estimated sources at two successive frequency bins, we test of the closeness of $\mathbf{G}(f_n)\mathbf{G}(f_{n-1})^{\#}$ to a diagonal matrix. Indeed, by using the representation (19), one can find the permutation matrix by minimizing :

$$\mathbf{\Pi}(f_n) = \arg \min_{\tilde{\mathbf{\Pi}}} \left\{ \sum_{i \neq j} \left(\tilde{\mathbf{\Pi}}\mathbf{C}(f_n)\mathbf{B}(f_n)\mathbf{G}(f_{n-1})^{\#} \right)_{ij}^2 \right\} . \quad (20)$$

In our simulations, we have used an exhaustive search to solve (20). However, when the number of sources is large, the exhaustive search becomes prohibitive. In that case, one can estimate $\mathbf{\Pi}(f_n)$ as the matrix with ones at the ij^{th} entry satisfying $|\mathcal{M}(f_n)|_{ij} = \max_k |\mathcal{M}(f_n)|_{ik}$ and zeros elsewhere with $\mathcal{M}(f_n) = \mathcal{C}(f_n)\mathbf{B}(f_n)\mathbf{G}(f_{n-1})^\#$. This solution has the advantage of simplicity but may lead to erroneous solution in difficult context. An alternative solution would be to decompose $\mathbf{\Pi}(f_n)$ as product of elementary permutations¹ $\mathbf{\Pi}_{(pq)}$. The latter is considered at a given iteration, only if it decrease criterion (20), if

$$|\mathcal{M}(f_n)|_{pq}^2 + |\mathcal{M}(f_n)|_{qp}^2 > |\mathcal{M}(f_n)|_{pp}^2 + |\mathcal{M}(f_n)|_{qq}^2$$

Finally, we obtain :

$$\mathbf{\Pi}(f_n) = \prod_{\text{nb of iterations}} \prod_{1 \leq p < q \leq N} \widetilde{\mathbf{\Pi}}_{(pq)} , \quad (21)$$

$\widetilde{\mathbf{\Pi}}_{(pq)}$ being either the identity matrix or the above permutation matrix $\mathbf{\Pi}_{(pq)}$ depending on the binary decision rule define above. We stop the iterative process, when all matrices $\widetilde{\mathbf{\Pi}}_{(pq)}$ are equal to the identity. We have observed that one or, at most, two iterations are sufficient to get the desired permutation. Finally, we apply the updated separation matrix $\mathbf{G}(f)$ to the frequency domain mixture :

$$\mathcal{S}_{\mathbf{s}}(t, f) = \mathbf{G}(f)\mathcal{S}_{\mathbf{x}}(t, f) . \quad (22)$$

4 Simulation results

We present here some numerical simulations to evaluate the performance of our algorithm. We consider an array of $M = 2$ sensors receiving two audio signals in the presence of stationary temporally white noise of covariance $\sigma^2\mathbf{I}$ (σ^2 being the noise power). 10000 samples are used with a sampling frequency of 8Khz (this represents 1.25sec recording). In order to evaluate the performance in the instantaneous mixture case, the separation quality is measured using the *Interference to Signal Ratio* (ISR) criterion [2] defined as :

$$ISR \stackrel{\text{def}}{=} \sum_{p \neq q} \frac{E(|(\mathbf{BA})_{pq}|^2) \rho_q}{E(|(\mathbf{BA})_{pp}|^2) \rho_p} \quad (23)$$

where $\rho_i = E(|s_i(t)|^2)$ is the i^{th} source power evaluated here as $\frac{1}{T} \sum_{t=0}^{T-1} |s_i(t)|^2$. Fig. 1-(a) represents the two original sources and their mixtures in the noiseless case. In Fig. 1-(b), we compare the performance of the proposed algorithm in instantaneous mixture case, to the Relative Newton algorithm developed by Zibulevsky et al. in [9] where the case of sparse sources is considered and to SOBI algorithm developed by Belouchrani et al. in [2]. We plot the residual interference between separated sources (ISR) versus the SNR. It is clearly shown that our algorithm (ISBS) performs better in terms of ISR especially for low SNRs as compared to the two other methods. In Fig. 2-(a), we represent the evolution of

¹ $\mathbf{\Pi}_{(pq)}$ is defined such as way that for a given vector \mathbf{y} , $\tilde{\mathbf{y}} = \mathbf{\Pi}_{(pq)}\mathbf{y}$ iff $\tilde{y}(k) = y(k)$, for $k \notin \{p, q\}$, $\tilde{y}(p) = y(q)$ and $\tilde{y}(q) = y(p)$.

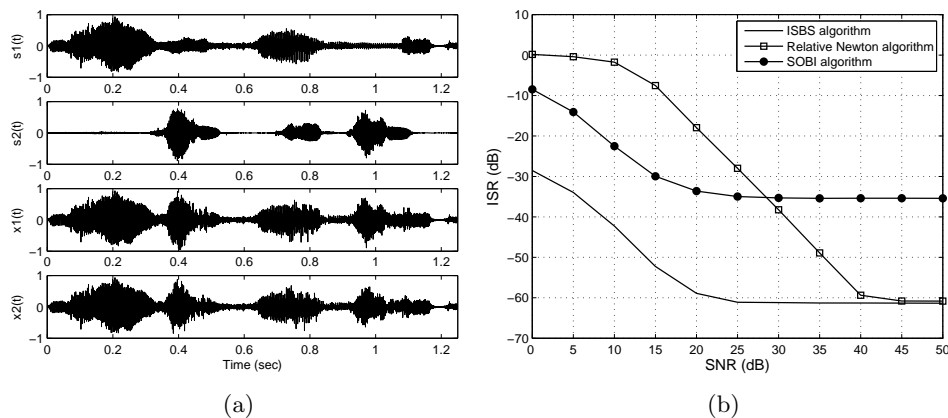


Fig. 1. (a) Up the two original source signals and bottom the two signal mixtures. (b) Interference to Signal Ratio (ISR) versus SNR for 2 audio sources and 2 sensors in instantaneous mixture case.

the ISR as a function of the iteration number. A fast convergence rate is observed. In Fig. 2-(b), we compare, in the 2×2 convolutive mixture case the separation performance of our algorithm, Deville’s algorithm in [13], Parra’s algorithm in [14] and extended version of Zibulevsky’s algorithm to the convolutive mixture case. The filter coefficients are chosen randomly and the channel order is $L = 128$. We use in this experiment the ISR criterion defined for the convolutive case in [14] that takes into account the fact the source estimates are obtained up to a scalar filter. We observe a significant performance gain in favor of the proposed method especially at low SNR values. Moreover, the complexity of the proposed algorithm is equal to $2N^2T + \mathcal{O}(N^2)$ flops per iteration whereas the complexity of the Relative Newton algorithm in [9] is $2N^4 + N^3T + N^6/6$.

5 Conclusion

This paper presents a blind source separation method for sparse sources in instantaneous mixture case and its extension to the convolutive mixture case. A sparse contrast function is introduced and an iterative algorithm based on gradient technique is proposed to minimize it and perform the BSS. Numerical simulation results have been given evidence the usefulness of the method. The proposed technique outperforms existing solutions in terms of separation quality and computational cost in both instantaneous and convolutive mixture cases.

References

1. Cardoso, J.F.: Blind signal separation : statistical principles. Proceedings of the IEEE **86**(10) (Oct. 1998) 2009–2025

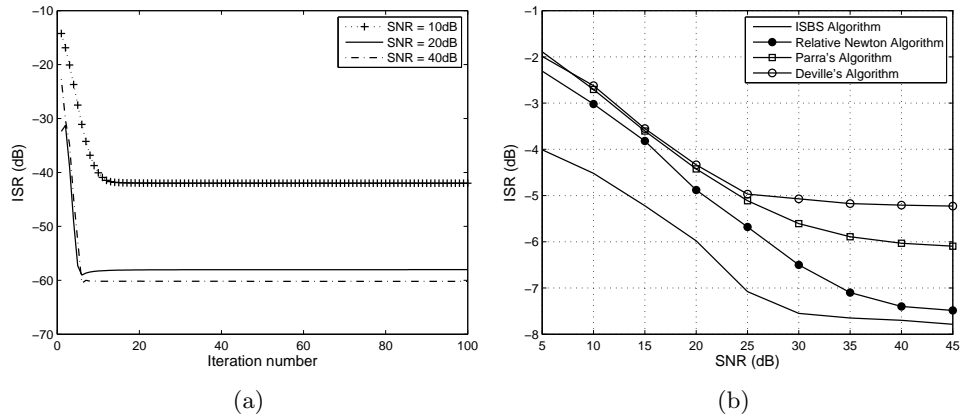


Fig. 2. (a) ISR as a function of the iteration number for 2 audio sources and 2 sensors in instantaneous mixture case. (b) ISR versus SNR for 2×2 convolutive mixture case.

2. Belouchrani, A., Abed-Meraim, K., Cardoso, J.F., Moulines, E.: A blind source separation technique using second-order statistics. *IEEE T-SP* **45**(2) (Feb. 1997)
3. Belouchrani, A., Amin, M.G.: Blind source separation based on time-frequency signal representations. *IEEE Transactions on Signal Processing* **46**(11) (Nov. 1998)
4. Abed-Meraim, K., Xiang, Y., Manton, J.H., Hua, Y.: Blind source separation using second order cyclostationary statistics. *IEEE T-SP* **49**(4) (April 2001) 694–701
5. Yilmaz, O., Rickard, S.: Blind separation of speech mixtures via time-frequency masking. *IEEE Transactions on Signal Processing* **52**(7) (July 2004) 1830–1847
6. Pham, D.T., Cardoso, J.F.: Blind separation of instantaneous mixtures of non stationary sources. *IEEE Transactions on Signal Processing* **49** (2001) 1837–1848
7. Smith, D., Lukasiak, J., Burnett, I.S.: An analysis of the limitations of blind signal separation application with speech. *Signal Processing* **86**(2) (Feb. 2006) 353–359
8. Cichocki, A., Amari, S.: Chapter 2. In: *Adaptive Blind Signal and Image Processing*. Wiley & Sons, Ltd., UK (2003)
9. Zibulevsky, M.: Sparse source separation with relative Newton method. In: *Proc. ICA*. (April 2003) 897–902
10. Pham, D.T., Garat, P.: Blind separation of mixture of independent sources through a quasi-maximum likelihood approach. *IEEE T-SP* **45**(7) (July 1997) 1712–1725
11. Murata, N., Ikeda, S., Ziehe, A.: An approach to blind source separation based on temporal structure of speech signals. *Neurocomputing* **41**(1-4) (Oct. 2001) 1–24
12. Pham, D.T., Servière, C., Boumaraf, H.: Blind separation of convolutive audio mixtures using nonstationarity. In: *Proc. ICA, Nara, Japan* (April 2003) 981–986
13. Albouy, B., Deville, Y.: Alternative structures and power spectrum criteria for blind segmentation and separation of convolutive speech mixtures. In: *Proc. ICA, Nara, Japan* (April 2003) 361–366
14. Parra, L., Spence, C.: Convolutive blind separation of non-stationary sources. *IEEE Transactions on Speech and Audio Processing* **8**(3) (May 2000) 320–327