

Adaptive blind separation of underdetermined mixtures based on sparse component analysis

YANG ZuYuan, HE ZhaoShui, XIE ShengLi[†] & FU YuLi

School of Electronics & Information Engineering, South China University of Technology, Guangzhou 510640, China

The independence *priori* is very often used in the conventional blind source separation (BSS). Naturally, independent component analysis (ICA) is also employed to perform BSS very often. However, ICA is difficult to use in some challenging cases, such as underdetermined BSS or blind separation of dependent sources. Recently, sparse component analysis (SCA) has attained much attention because it is theoretically available for underdetermined BSS and even for blind dependent source separation sometimes. However, SCA has not been developed very sufficiently. Up to now, there are only few existing algorithms and they are also not perfect as well in practice. For example, although Lewicki-Sejnowski's natural gradient for SCA is superior to K-mean clustering, it is just an approximation without rigorously theoretical basis. To overcome these problems, a new natural gradient formula is proposed in this paper. This formula is derived directly from the cost function of SCA through matrix theory. Mathematically, it is more rigorous. In addition, a new and robust adaptive BSS algorithm is developed based on the new natural gradient. Simulations illustrate that this natural gradient formula is more robust and reliable than Lewicki-Sejnowski's gradient.

underdetermined mixtures, blind source separation (BSS), dependent sources, sparse component analysis (SCA), sparse representation, independent component analysis (ICA), natural gradient

1 Introduction

Since blind source separation (BSS) has many potential applications in digital communication, speech and image processing, array signal processing, and medical signal processing, etc. It has become a hot spot in both signal processing and neural network fields. The mathematical model of BSS can be formulated as follows:

$$\mathbf{x} = \mathbf{A}\mathbf{s}, \quad (1)$$

Received November 11, 2006; accepted September 11, 2007

doi: 10.1007/s11432-008-0030-4

[†]Corresponding author (email: adshlxie@scut.edu.cn)

Supported by the National Natural Science Foundation of China (Grant Nos. 60505005, 60674033, 60774094 and U0635001), the Natural Science Fund of Guangdong Province, China (Grant Nos. 05103553 and 05006508), the Postdoctoral Science Foundation for Innovation from South China University of Technology, and China Postdoctoral Science Foundation (Grant No. 20070410237)

where $\mathbf{x} = (x_1, \dots, x_m)^T$ is the vector of observed signals, $\mathbf{s} = (s_1, \dots, s_n)^T$ is the vector of unknown sources, and \mathbf{A} is the unknown $m \times n$ mixing matrix. The objective of BSS is to estimate source signals up to a scaling factor and a permutation factor only from the mixing signals.

The conventional BSS is usually based on the independence assumption^[1,2]. For example, independent component analysis (ICA) based BSS^[3-5] assumes that the sources are statistically mutually independent and requires that the number of sensors is not less than the number of sources, i.e., $m \geq n$. For nonsingular square matrix \mathbf{A} , the problem corresponds to estimating the inverse matrix \mathbf{A}^{-1} . Under this framework, many algorithms were proposed, e.g., information theory based algorithms (Bell-Sejnowski^[6], Xie^[7]), higher order statistics based algorithms (Hyvarinen^[8-11], Comon^[12]), non-stationarity based algorithms (Hyvarinen^[12], Matsuoka^[13]), and the geometrical algorithms^[14]. However, for the underdetermined BSS or ill-posed BSS ($m < n$), the inverse matrix \mathbf{A}^{-1} does not exist. The above methods are unavailable. Even if the mixing matrix \mathbf{A} is given, one cannot obtain the unique solution in this case. Thus, new methods should be developed for estimating \mathbf{A} and \mathbf{s} .

Furthermore, since the number of source signals is unknown and we usually have an only finite number of observed signals, sometimes the underdetermined mixtures occur in practice. Therefore, it is necessary to study the underdetermined models. Blind extraction was first proposed to perform underdetermined BSS, which extracts the source signals one by one. It should be noted that at most a part of sources, but not all, can be estimated by blind extraction for underdetermined BSS^[15,16]. On the other hand, the precondition "statistical independence assumption" is not always reasonable in practice. For example, EEG, MEG, and FMRI sources are not statistically independent very often. BSS methods based on ICA are not available for these challenging problems because of the dependent sources.

For these reasons, the novel tools beyond ICA should be considered. Among them, the sparsity of source signals was recently exploited to achieve BSS by sparse representation or sparse decomposition. In this way, sparse component analysis (SCA) was directly or indirectly applied to underdetermined BSS, and promising results were reported^[17]. For example, Belouchrani^[18] proposed a maximum likelihood method, Zibulevsky^[17] proposed a sparse decomposition method, overcomplete representation method was proposed by Lee^[19], Lewicki^[20], and Li^[21], and the sparse representation in frequency domain was proposed by Bofill and Zebulevsky^[8]. Typically, six source signals were successfully separated from only two mixing signals by Bofill and Zebulevsky^[22]. In addition, SCA is even available for blind separation of dependent source signals sometimes (see Simulation 1 in this paper)^[17].

As a whole, there are mainly two classes of BSS algorithms based on SCA. One of them is the so-called two-stage clustering-then- l^1 -optimization algorithm, the other is adaptive algorithm based on overcomplete sparse representation. As its name implies, the former has two steps^[17,21]: clustering analysis gives the estimation of the channel and linear programming estimates the source signals after the mixing matrix is estimated. Generally speaking, many clustering algorithms including K-mean are not very accurate and robust in practice. Thus, Lee, Lewicki et al.^[17,23] developed the overcomplete representation based adaptive method, which can fill the gap of the former. Moreover, the results by Zhang, Amari, and Cichocki^[24] further strengthen the adaptive method based on Lewicki-Sejnowski natural gradient^[20]. However, Lewicki-Sejnowski natural gradient was derived by repeated approximations and is not rigorous mathematically.

Start from the cost function of SCA directly, we develop a rigorous natural gradient formula based on matrix theory. Also, an SCA based algorithm is presented for the underdetermined BSS in this paper. Simulations illustrate availability of this algorithm. Comparing with Lewicki-Sejnowski's natural gradient, the proposed natural gradient is more robust.

2 Sparse component analysis (SCA)

The SCA (or sparse representation) model is the same as that of BSS (1). However, different from ICA, SCA aims at estimating the basis matrix \mathbf{A} and sparse components \mathbf{s} such that \mathbf{s} is as sparse as possible. In other words, \mathbf{s} takes as many zeros as possible, or \mathbf{s} has few nonzero entries if possible. Mathematically, this problem can be expressed as an l^0 -norm optimization one^[21,25],

$$(P_0) \min J(\mathbf{A}, \mathbf{s}) = \min_{\mathbf{A}, \mathbf{s}} \|\mathbf{s}\|^0 = \min_{\mathbf{A}, \mathbf{s}} \sum_{i=1}^n |s_i|^0, \quad \text{st: } \mathbf{A}\mathbf{s} = \mathbf{x}, \quad (2)$$

where $\|\mathbf{s}\|^0$ measures the number of nonzero entries in the vector \mathbf{s} ^[21,25]. Here $|s_i|^0 = 1$ for $s_i \neq 0$, while $|s_i|^0 = 0$ if and only if $s_i = 0$.

Although the problem (2) can give the sparsest solution, it has no unique solution in general. Furthermore, it is very difficult to solve problem (2), and its solution is not robust but is very sensitive to noise. For this reason, Donoho^[25] and Li^[21] suggested using l^1 -norm to measure the sparsity. Thus, the l^0 -norm optimization problem (2) is replaced by its approximation—an l^1 -norm optimization problem

$$(P_1) \min J(\mathbf{A}, \mathbf{s}) = \min_{\mathbf{A}, \mathbf{s}} \|\mathbf{s}\|^1 = \min_{\mathbf{A}, \mathbf{s}} \sum_{i=1}^n |s_i|, \quad \text{st: } \mathbf{A}\mathbf{s} = \mathbf{x}. \quad (3)$$

Donoho and Li found that the solutions to (3) are good approximations of the solutions to (2), when the source \mathbf{s} is sufficiently sparse. Furthermore, there are many existing mathematical tools for l^1 -norm optimization problem. So we are interested to achieve SCA by solving problem (3). An example was given by Zibulevsky and Pearlmutter^[17].

In the engineering applications, many signals (e.g., speech signals or other signals detected by some instruments) are with limited frequency bandwidths and are a little smooth in some degree (but maybe are not very smooth), which results in that these signals (except white Gaussian noise) in the real world are sparse in some degree in the time domain, or they can be sparsely represented in an appropriate transform domain (e.g., frequency domain or wavelet domain). Thus, SCA has many applications in practice^[15,17]. Especially, it can be applied to BSS. If the source signals are sufficiently sparse, SCA can be applied to underdetermined BSS or BSS of dependent sources (see simulation 1 in this paper)^[15,17,19,21,22]. SCA can achieve better result than the conventional ICA even for the well-determined BSS^[17,19]. To well satisfy the sparseness assumption, if the source signals are not sparse in time domain, we can consider to perform BSS in an appropriate transform domain by Fourier transform, wavelet or other transforms.

Next, we discuss how to perform SCA by minimizing l^1 -norm, and here we mainly discuss how solve problem (3).

3 Natural gradient of mixing matrix \mathbf{A}

To identify mixing matrix \mathbf{A} , Lewicki and Sejnowski^[20] originally presented an approximate natural gradient with respect to the mixing matrix \mathbf{A} based on “overcomplete representation

learning" theory. Later, Lee and Lewicki^[19] successfully identified A using this approximate gradient. Here start from the optimization problem (3), we will derive a new and rigorous natural gradient for learning the mixing matrix A . For this problem, we have the following theorem:

Theorem 1. For the optimization problem (3), the natural gradient of cost function $J(\cdot)$ with respect to A can be expressed as follows:

$$\nabla J(A) = -A \cdot (\text{sign}(s) \cdot s^T), \quad (4)$$

where $\text{sign}(s) = (\text{sign}(s_1), \dots, \text{sign}(s_n))^T$, and $\text{sign}(\cdot)$ is a sign function.

Proof. In the optimization problem (3), the cost function $J = \sum_{i=1}^n |s_i|$ is not differentiable at point 0. Thus, we introduce an auxiliary function $\rho(x, \delta)$

$$\rho(x, \delta) = \begin{cases} |x|, & |x| \geq \delta, \\ \frac{x^2}{2\delta} + \frac{\delta}{2}, & |x| \leq \delta, \end{cases} \quad (5)$$

where δ is a small positive number. The function $\rho(x, 0.1)$ is shown in Figure 1.

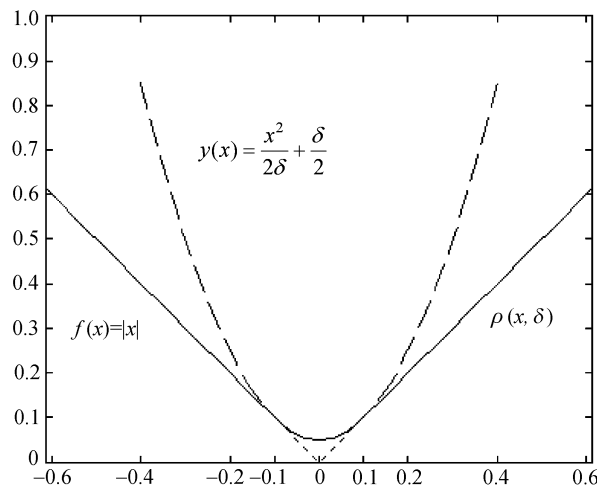


Figure 1 Function $\rho(x, \delta)$ when $\delta = 0.1$.

Function $\rho(x, \delta)$ has the following properties:

- 1) $\rho(x, \delta) \geq |x|$, for $\delta \rightarrow 0$; $\rho(x, \delta)$ uniformly converges to $|x|$, i.e., $\lim_{\delta \rightarrow 0} \rho(x, \delta) = |x|$, $x \in (-\infty, +\infty)$;
- 2) $\rho(x, \delta)$ is differentiable with respect to x , and its derivative is as follows:

$$\phi(x, \delta) = \frac{d\rho}{dx} = \begin{cases} \text{sign}(x), & |x| \geq \delta, \\ \frac{x}{\delta}, & |x| \leq \delta, \end{cases} \quad (6)$$

where $\phi(0, \delta) = 0$ and $\phi(x, \delta)$ converges to $\text{sign}(x)$ almost everywhere, i.e., $\forall x \in (-\infty, +\infty)$, $\lim_{\delta \rightarrow 0} \phi(x, \delta) = \text{sign}(x)$.

By the function $\rho(x, \delta)$, problem (3) can be approximated by the following optimization one:

$$\min_{A,s} L(A,s) = \min_{A,s} \sum_{i=1}^n \rho(s_i, \delta), \quad \text{st: } As = x, \quad (7)$$

when $\delta \rightarrow 0$, problem (7) tends to problem (3). Since function $\rho(x, \delta)$ is differentiable with respect to x , we can compute the derivative of $L(\cdot)$ with respect to x .

We prove Theorem 1 for two cases: $m=n$ (A is nonsingular), and $m < n$ (A is overcomplete).

Case 1: $m=n$ and A is nonsingular

Since A is nonsingular, suppose that $A = W^{-1}$. From (1), we have

$$s = Wx. \quad (8)$$

Then, the optimization problem (7) can be changed into

$$\min_{W} L(W, \delta) = \min_{W} \sum_{i=1}^n \rho(s_i, \delta), \quad \text{st: } s = Wx. \quad (9)$$

The derivative of $L(\cdot, \delta)$ with respect to w_{ij} is

$$\frac{\partial L(\cdot, \delta)}{\partial w_{ij}} = \sum_{k=1}^n \frac{\partial L(\cdot, \delta)}{\partial s_k} \cdot \frac{\partial s_k}{\partial w_{ij}} = \frac{\partial L(\cdot, \delta)}{\partial s_i} \cdot \frac{\partial s_i}{\partial w_{ij}} = \phi(s_i, \delta) \cdot x_j, \quad i, j = 1, \dots, n.$$

That is,

$$\frac{\partial L(\cdot, \delta)}{\partial W} = \phi(s, \delta) \cdot x^T, \quad (10)$$

where $\phi(s, \delta) = (\phi(s_1, \delta), \dots, \phi(s_n, \delta))^T$. The function $\phi(\cdot, \delta)$ is defined as (6). Since $A = W^{-1}$, the gradient of $L(\cdot, \delta)$ with respect to A can be written as (see Appendix),

$$\frac{\partial L(\cdot, \delta)}{\partial A} = -A^{-T} \cdot \frac{\partial L(\cdot, \delta)}{\partial W} \cdot A^{-T} = -A^{-T} \cdot \phi(s, \delta) \cdot x^T W^T = -A^{-T} \cdot \phi(s, \delta) \cdot s^T.$$

Thus, the natural gradient of $L(\cdot, \delta)$ with respect to A is

$$\nabla L(A, \delta) = AA^T \cdot \frac{\partial L(\cdot, \delta)}{\partial A} = -AA^T \cdot A^{-T} \cdot \phi(s, \delta) \cdot s^T = -A \cdot \phi(s, \delta) \cdot s^T. \quad (11)$$

Note that there is no inverse matrix A^{-1} in the natural gradient $\nabla L(A, \delta)$ in (11). So the natural gradient $\nabla L(A, \delta)$ is simpler than standard gradient $\partial L(\cdot, \delta) / \partial A$. From the Lemma 2 in Appendix, this natural gradient can be used to solve optimization problems (7). When $\delta \rightarrow 0$, problem (7) approximates to problem (3), $L(\cdot, \delta) \rightarrow J(\cdot)$ and $\nabla L(A, \delta)$ tends to $\nabla J(A)$ almost everywhere, i.e.,

$$\forall s, \lim_{\delta \rightarrow 0} \nabla L(A, \delta) = \nabla J(A) = -A \cdot (\text{sign}(s) \cdot s^T). \quad (12)$$

Thus, the natural gradient of $J(\cdot)$ with respect to A is $\nabla J(A) = -A \cdot (\text{sign}(s) \cdot s^T)$. Theorem 1 is proven in Case 1.

Case 2: $m < n$

In this case, A is an $m \times n$ matrix. We add $n - m$ virtual observed signals $x_{(n-m) \times 1}^{\text{virtual}}$. The model (1) can be extended as

$$x^{\text{new}} = A^{\text{new}} s. \quad (13)$$

It is possible to add the virtual observed signals such that A^{new} is an $n \times n$ nonsingular matrix,

where $\mathbf{x}^{\text{new}} = \begin{pmatrix} \mathbf{x} \\ \mathbf{x}_{(n-m) \times 1}^{\text{virtual}} \end{pmatrix}$, and $\mathbf{A}^{\text{new}} = \begin{pmatrix} \mathbf{A} \\ \bar{\mathbf{A}}_{(n-m) \times n}^{\text{virtual}} \end{pmatrix}$. The sub-matrix $\bar{\mathbf{A}}_{(n-m) \times n}^{\text{virtual}}$ associates with $\mathbf{x}_{(n-m) \times 1}^{\text{virtual}}$. Since \mathbf{A}^{new} is nonsingular, from the proof of Case 1, the natural gradient of $J(\cdot)$ with respect to \mathbf{A}^{new} is analogously as follows:

$$\begin{aligned} \nabla J(\mathbf{A}^{\text{new}}) &= \begin{pmatrix} \nabla J(\mathbf{A}) \\ \nabla J(\bar{\mathbf{A}}_{(n-m) \times n}^{\text{virtual}}) \end{pmatrix} = -\mathbf{A}^{\text{new}} \cdot \text{sign}(\mathbf{s}) \cdot \mathbf{s}^T \\ &= -\begin{pmatrix} \mathbf{A} \\ \bar{\mathbf{A}}_{(n-m) \times n}^{\text{virtual}} \end{pmatrix} \cdot \text{sign}(\mathbf{s}) \cdot \mathbf{s}^T = \begin{pmatrix} -\mathbf{A} \cdot \text{sign}(\mathbf{s}) \cdot \mathbf{s}^T \\ -\bar{\mathbf{A}}_{(n-m) \times n}^{\text{virtual}} \cdot \text{sign}(\mathbf{s}) \cdot \mathbf{s}^T \end{pmatrix}, \end{aligned}$$

or, it can be re-written as

$$\nabla J(\mathbf{A}^{\text{new}}) = \begin{pmatrix} \nabla J(\mathbf{A}) \\ \nabla J(\bar{\mathbf{A}}_{(n-m) \times n}^{\text{virtual}}) \end{pmatrix} = \begin{pmatrix} -\mathbf{A} \cdot \text{sign}(\mathbf{s}) \cdot \mathbf{s}^T \\ -\bar{\mathbf{A}}_{(n-m) \times n}^{\text{virtual}} \cdot \text{sign}(\mathbf{s}) \cdot \mathbf{s}^T \end{pmatrix}.$$

Thus, for $m < n$, the natural gradient of $J(\cdot)$ with respect to \mathbf{A} also can be formulated as

$$\nabla J(\mathbf{A}) = -\mathbf{A} \cdot (\text{sign}(\mathbf{s}) \cdot \mathbf{s}^T).$$

The natural gradient (4) is very simple and effective, which can be employed to perform SCA. In the next section, we will establish a BSS algorithm based on SCA.

4 BSS algorithm

By the above natural gradient, we have the following iterative rule for learning the mixing matrix \mathbf{A} :

$$\mathbf{A}^{(k+1)} = \mathbf{A}^{(k)} - \mu \nabla J(\mathbf{A}^{(k)}) = \mathbf{A}^{(k)} + \mu \mathbf{A}^{(k)} \cdot \text{sign}(\mathbf{s}^{(k)}) \cdot (\mathbf{s}^{(k)})^T, \quad (14)$$

where $\mu > 0$. As pointed out in ref. [17], the solution to problem (3) will be $\mathbf{s} \approx 0$, if there are no additional constraints. Obviously, this is not the right solution. For this point, we add some additional constraints to problem (3)^[21,26]. We constrain the ℓ^2 -norm of each column of \mathbf{A} to be 1^[26], i.e., $\|\mathbf{a}_i\| = 1, i = 1, \dots, n$. To this end, we consider a similar method to Parra and Spence's^[27] constrained gradient. In formula (14), instead of natural gradient $\nabla J(\mathbf{A})$, we consider the projection of $\nabla J(\mathbf{A})$ on the hyperplanes defined by $\|\mathbf{a}_i\| = 1, i = 1, \dots, n$. The projection operator for the i th column $\nabla J(\mathbf{a}_i)$ of $\nabla J(\mathbf{A}) = [\nabla J(\mathbf{a}_1), \dots, \nabla J(\mathbf{a}_n)]$ is

$$P_i = \mathbf{I} - \mathbf{a}_i \cdot \mathbf{a}_i^T, \text{ or } P: \nabla J(\mathbf{a}_i) \rightarrow (\mathbf{I} - \mathbf{a}_i \cdot \mathbf{a}_i^T) \cdot \nabla J(\mathbf{a}_i). \quad (15)$$

We can also re-write it as a constrained natural gradient

$$\nabla J(\mathbf{A}) \Big|_{\|\mathbf{a}_i\|=1, i=1, \dots, n} = \nabla J(\mathbf{A}) - \mathbf{A} \cdot \text{diag}(\mathbf{A}^T \cdot \nabla J(\mathbf{A})). \quad (16)$$

The constrained natural gradient (16) $\nabla J(\mathbf{A}) \Big|_{\|\mathbf{a}_i\|=1, i=1, \dots, n}$ is not always necessary. From numerical experiments, we find that good performance also can be achieved generally if we combine natural gradient $\nabla J(\mathbf{A})$ with the column normalization operators.

In formula (14), we cannot update $\mathbf{A}^{(k)}$ unless $\mathbf{s}^{(k)}$ is given. To achieve this goal, we can

alternatively update \mathbf{A} and \mathbf{s} in a similar way in ref. [19]. If the mixing matrix \mathbf{A} is known, the optimization problem (3) can be simplified as

$$\begin{cases} \min J(\mathbf{s} | \mathbf{A}) = \min_{\mathbf{s}} \sum_{i=1}^n |s_i|, \\ \text{st} : \mathbf{A}\mathbf{s} = \mathbf{x}. \end{cases} \quad (17)$$

The shortest path decomposition (SPD) can be employed to solve the optimization problem (17)^[22,26]. For simplicity, we denote the solution of (17) as $\hat{\mathbf{s}}^* = \text{SPD}(\mathbf{A}, \mathbf{x})$.

From the above discussion, the BSS algorithm for underdetermined mixtures based on SCA can be outlined as follows:

- 1) Initialize \mathbf{A} as $\mathbf{A}^{(0)}$, where each column of $\mathbf{A}^{(0)}$ is normalized such that its ℓ^2 -norm is "1". Set the step size $\mu > 0$ and let $k = 0$.
- 2) Estimate $\hat{\mathbf{s}}^{(k)} = \text{SPD}(\mathbf{A}^{(k)}, \mathbf{x})$ by SPD.
- 3) Substituting $\hat{\mathbf{s}}^{(k)}$ into (14), compute the constrained natural gradient $\nabla J(\mathbf{A})|_{\|\mathbf{a}_i\|=1, i=1, \dots, n}$. Update \mathbf{A} as $\mathbf{A}^{(k+1)}$. Let $k = k + 1$.
- 4) If the algorithm is convergent, turn to 5), otherwise, go to 2).
- 5) Output the estimation of the source signals $\mathbf{s}^* = \hat{\mathbf{s}}^{(k)}$.

5 Simulations

To evaluate performance, two performance indices are used here. One is the bias angle $ae(\mathbf{a}, \hat{\mathbf{a}})$ between the column \mathbf{a} in mixing matrix and its estimation $\hat{\mathbf{a}}$. It is defined as

$$ae(\mathbf{a}, \hat{\mathbf{a}}) = \frac{180}{\pi} \arccos \left(\frac{\langle \mathbf{a}, \hat{\mathbf{a}} \rangle}{\|\mathbf{a}\| \cdot \|\hat{\mathbf{a}}\|} \right), \quad (18)$$

$ae(\mathbf{a}, \hat{\mathbf{a}})$ is small if the direction of \mathbf{a} is close to that of $\hat{\mathbf{a}}$. Otherwise, $ae(\mathbf{a}, \hat{\mathbf{a}})$ is large.

Another index is signal to interference ratio (SIR) between separated signal $\hat{\mathbf{s}}$ and source signal \mathbf{s} ^[22].

$$\text{SIR}(\mathbf{s}, \hat{\mathbf{s}}) = 10 \log \left[\frac{s^2}{(s - \hat{\mathbf{s}})^2} \right] \text{ (dB)}. \quad (19)$$

Since separated signal $\hat{\mathbf{s}}$ can be different from \mathbf{s} up to a scaling factor c ($c \neq 0$), as treated in ref. [22], we adjust \mathbf{s} and $\hat{\mathbf{s}}$ to be with the equal energy, before the SIR is computed.

Simulation 1. SCA based BSS for dependent signals

Consider two source signals $s_1(t)$ and $s_2(t)$, $t = 1, \dots, 1000$. All the samples of $s_1(t)$ and $s_2(t)$ take values from the binary set $\{0, 1\}$, i.e., $s_1(t) = 1$ or $s_1(t) = 0$, and $s_2(t) = 1$ or $s_2(t) = 0$. In addition, we constrain that $s_1(t) + s_2(t) = 1$, $t = 1, \dots, 1000$. Obviously, $s_1(t)$ and $s_2(t)$ are statistically dependent because they satisfy the equation $s_1(t) = 1 - s_2(t)$, $t = 1, \dots, 1000$. The mixing matrix is generated by random as

$$\mathbf{A} = \begin{pmatrix} 0.9673 & -0.9105 \\ 0.2537 & 0.4135 \end{pmatrix}.$$

From model (1), we have observed signals $x_1(t)$ and $x_2(t)$.

Set the step $\mu = 0.01$. The algorithm converged after 500 iterations. The mixing matrix \mathbf{A} is estimated as follows:

$$\hat{\mathbf{A}} = \begin{pmatrix} 0.9671 & -0.9093 \\ 0.2543 & 0.4160 \end{pmatrix}.$$

By calculation, we have $ae(\mathbf{a}_1, \hat{\mathbf{a}}_1) = 0.036$ and $ae(\mathbf{a}_2, \hat{\mathbf{a}}_2) = 0.16$. The SIRs of the signals are 151.7109 dB and 109.8398 dB. Thus, two dependent signals were successfully separated.

Simulation 2. The performance comparison between the constrained natural gradient (16) and the natural gradient presented by Lewicki-Sejnowski.

Suppose that the signals follow Laplacian distribution. From ref. [19], Lewicki-Sejnowski's natural gradient is $\Delta \mathbf{A} \propto -\mathbf{A} \cdot (-\text{sign}(\mathbf{s}) \cdot \mathbf{s}^T + \mathbf{I})$, in which the identity matrix \mathbf{I} prevents the separated signals tending to zero. The projection gradient presented in this paper plays the same role.

For the BSS example of “two observed signals and three source signals” in ref. [13], simulations show that both our natural gradient and Lewicki-Sejnowski's natural gradient exhibited good performance. To further compare their performance, Bofill and Zibulevsky's^[22] example of “two observed signals mixed by six source signals” is taken here.

The sources are six flute signals. There are a total of 32768 samples in each signal. Figure 2(a) shows the waveforms of them.

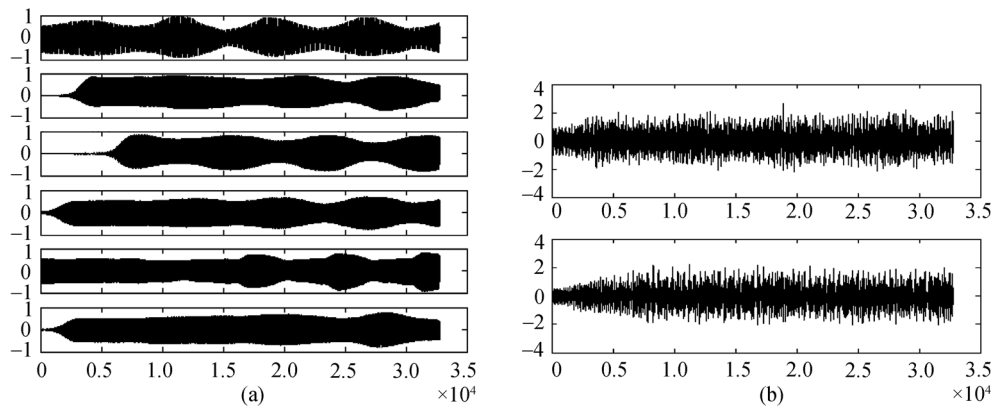


Figure 2 (a) The six source signals; (b) the two mixed signals.

Taking the same mixing matrix as in ref. [8],

$$\mathbf{A} = \begin{pmatrix} 0.9659 & 0.7071 & 0.2588 & -0.2588 & -0.7071 & -0.9659 \\ 0.2588 & 0.7071 & 0.9659 & 0.9659 & 0.7071 & 0.2588 \end{pmatrix},$$

the mixing signals are shown in Figure 2(b). For the two natural gradients, the same initialization was taken to update \mathbf{A} . Extensive tests show that Lewicki-Sejnowski's natural gradient failed to achieve BSS. For example, if the initialization is taken randomly as follows:

$$\mathbf{A}^{(0)} = \begin{pmatrix} -0.2514 & 0.3994 & -0.6935 & 0.9995 & 0.8823 & -0.2491 \\ -0.9679 & 0.9168 & 0.7204 & -0.0316 & 0.4708 & 0.9685 \end{pmatrix},$$

and the step size is set as $\mu = 0.01$, the algorithm that uses gradient (16) is convergent after 100

iterations. The algorithm by Lewicki-Sejnowski's natural gradient also converges after more than 200 iterations. Let the estimations of \mathbf{A} be $\hat{\mathbf{A}}_{\text{Our}}$ and $\hat{\mathbf{A}}_{\text{LS}}$, respectively. We have

$$\hat{\mathbf{A}}_{\text{Our}} = \begin{pmatrix} 0.9667 & 0.7096 & 0.2579 & -0.2597 & -0.7055 & -0.9657 \\ 0.2560 & 0.7046 & 0.9662 & 0.9657 & 0.7087 & 0.2595 \end{pmatrix},$$

$$\hat{\mathbf{A}}_{\text{LS}} = \begin{pmatrix} 6.1008 & 3.2397 & 0.0000 & -1.3476 & -4.5548 & -0.0620 \\ 1.4805 & 4.2478 & 0.0000 & 5.4940 & 3.8266 & 0.0165 \end{pmatrix}.$$

In this case, the bias angles and SIRs are given in Tables 1 and 2.

Table 1 The bias angles of each pair of columns of \mathbf{A} and $\hat{\mathbf{A}}$

Simulation 2	The angles (°)					
	$ae(a_1, \hat{a}_1)$	$ae(a_2, \hat{a}_2)$	$ae(a_3, \hat{a}_3)$	$ae(a_4, \hat{a}_4)$	$ae(a_5, \hat{a}_5)$	$ae(a_6, \hat{a}_6)$
Lewicki-Sejnowski's natural gradient	1.3597	7.6681	10.1678	1.2179	4.9657	0.1005
Constrained natural gradient (16)	0.1663	0.2024	0.0570	0.0521	0.1311	0.0399

Table 2 The SIRs of the separated signals

Simulation 2	SIR (dB)					
	1	2	3	4	5	6
Lewicki-Sejnowski's natural gradient	12.8135	9.0512	0	14.3304	10.4519	0
Constrained natural gradient (16)	50.5133	52.1840	49.1864	43.4974	49.1196	51.9210

From Tables 1, 2, and Figure 3, we can see that Lewicki-Sejnowski's natural gradient failed to achieve BSS, while the constrained natural gradient (16) succeeded.

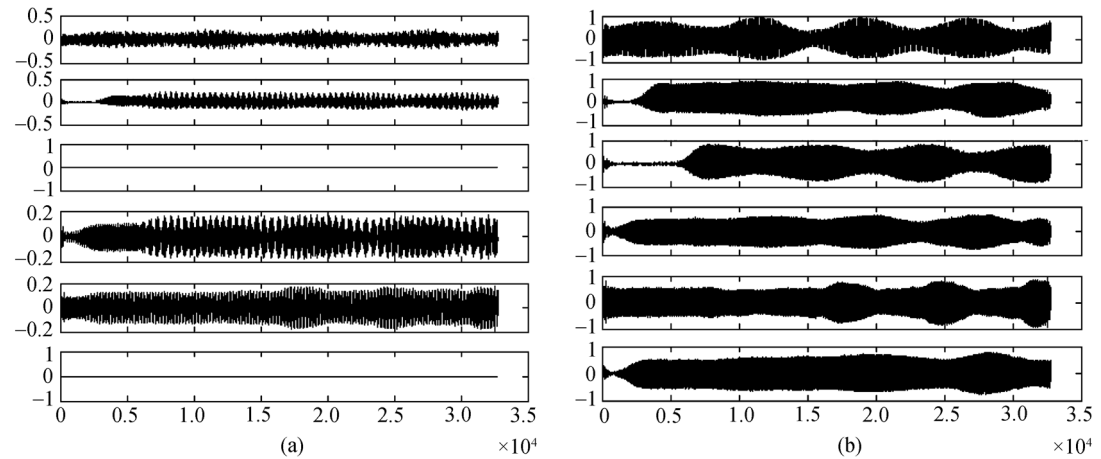


Figure 3 The separated signals by two natural gradients. (a) The separated signals by Lewicki-Sejnowski's gradient; (b) the separated signals by the constrained natural gradient (16).

Simulation 3. The robustness test of algorithm in noisy environment

The same source signals and mixing matrix are taken as in Simulation 2. However, 10 dbw white Gaussian noise was added to the observed signals here. We check the robustness of our algorithm in the noisy environment. The observed signals corrupted by noise are shown in Figure 4.

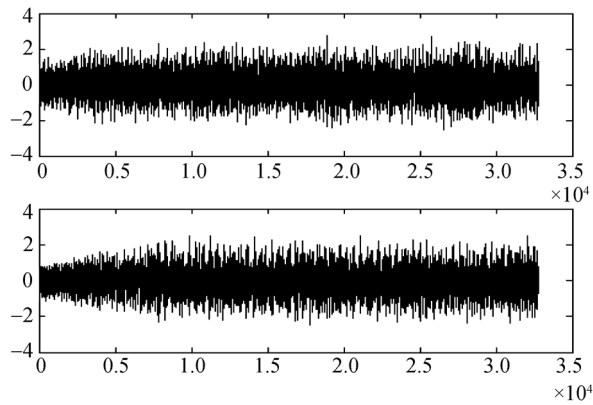


Figure 4 The observed signals with noise.

The learning step is set as $\mu = 0.01$. After 50 iterations, the algorithm converged and obtained the estimation of mixing matrix \mathbf{A} as follows:

$$\hat{\mathbf{A}} = \begin{pmatrix} 0.2557 & 0.7069 & -0.7093 & -0.9665 & 0.9670 & -0.2612 \\ 0.9668 & 0.7073 & 0.7049 & 0.2567 & 0.2549 & 0.9653 \end{pmatrix}.$$

The bias angles and SIRs are shown in Tables 3 and 4, respectively. Figure 5 shows the separated signals in noisy case.

Table 3 The bias angles of each pair of columns of \mathbf{A} and $\hat{\mathbf{A}}$ in noisy case

Bias angles $ae(\mathbf{a}, \hat{\mathbf{a}})$ ($^{\circ}$)					
$ae(\mathbf{a}_1, \hat{\mathbf{a}}_1)$	$ae(\mathbf{a}_2, \hat{\mathbf{a}}_2)$	$ae(\mathbf{a}_3, \hat{\mathbf{a}}_3)$	$ae(\mathbf{a}_4, \hat{\mathbf{a}}_4)$	$ae(\mathbf{a}_5, \hat{\mathbf{a}}_5)$	$ae(\mathbf{a}_6, \hat{\mathbf{a}}_6)$
0.2333	0.0144	0.1842	0.1393	0.1771	0.1282

Table 4 The SIRs of the estimations by the proposed natural gradient in noisy case

Simulation 3	SIRs (dB)					
Source signals	1	2	3	4	5	6
Estimated signals	24.8822	25.0302	25.1542	24.2955	24.8437	25.0588

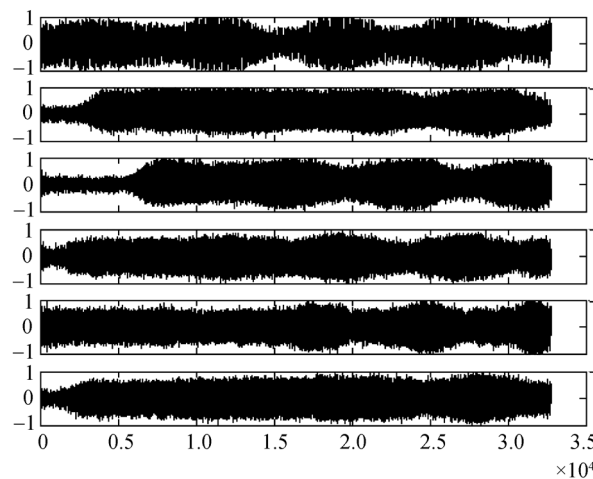


Figure 5 The separated signals (in noisy case).

Table 3 shows that although the heavy noise was added to the mixtures (10 dbw), the mixing matrix was very precisely estimated. Table 4 shows the SIRs of the separated signals. All SIRs are greater than 24 dB.

6 Conclusion

Most of BSS methods are based on the conventional ICA, which are unavailable for some challenging problems such as underdetermined BSS or blind separation of dependent sources. To face these challenges, SCA method is introduced in this paper. Also, this paper presents a new natural gradient and rigorously proves it by matrix theory. In addition, this natural gradient is successfully applied to SCA. The developed SCA can be used to perform underdetermined BSS if the sources are sufficiently sparse. Furthermore, our method is a potential tool for blind separation of dependent sources.

Appendix

Lemma 1. Let \mathbf{A} be an $n \times n$ invertible matrix and \mathbf{W} be the inverse of \mathbf{A} , i.e., $\mathbf{W} = \mathbf{A}^{-1}$. If $J(\mathbf{A})$ is a differentiable with respect to \mathbf{A} (with $n \times n$ variables), then the partial derivatives of J with respect to matrix \mathbf{A} and the inverse \mathbf{W} satisfy the following relation:

$$\frac{\partial J}{\partial \mathbf{A}} = -\mathbf{A}^{-T} \frac{\partial J}{\partial \mathbf{W}} \mathbf{A}^{-T}.$$

Proof. Obviously, we have

$$\mathbf{A}^T \mathbf{W}^T = \mathbf{I}_{n \times n}. \quad (\text{A1})$$

Let a_{ij} be the i th row and j th column of \mathbf{A} . We have

$$\begin{aligned} \frac{\partial \mathbf{A}^T}{\partial a_{ij}} \mathbf{W}^T + \mathbf{A}^T \frac{\partial \mathbf{W}^T}{\partial a_{ij}} &= \mathbf{0}_{n \times n} \Rightarrow \\ \frac{\partial \mathbf{W}^T}{\partial a_{ij}} &= -\mathbf{A}^{-T} \frac{\partial \mathbf{A}^T}{\partial a_{ij}} \mathbf{W}^T = -\mathbf{A}^{-T} \frac{\partial \mathbf{A}^T}{\partial a_{ij}} \mathbf{A}^{-T}. \end{aligned} \quad (\text{A2})$$

On the other hand, we have

$$\begin{aligned} \frac{\partial J}{\partial a_{ij}} &= \text{Trace} \left\{ \frac{\partial J}{\partial \mathbf{W}} \frac{\partial \mathbf{W}^T}{\partial a_{ij}} \right\} \\ &= -\text{Trace} \left\{ \left(\frac{\partial J}{\partial \mathbf{W}} \mathbf{A}^{-T} \frac{\partial \mathbf{A}^T}{\partial a_{ij}} \right) \mathbf{A}^{-T} \right\} = -\text{Trace} \left\{ \left(\mathbf{A}^{-T} \frac{\partial J}{\partial \mathbf{W}} \mathbf{A}^{-T} \right) \frac{\partial \mathbf{A}^T}{\partial a_{ij}} \right\}, \end{aligned} \quad (\text{A3})$$

where $\frac{\partial \mathbf{A}^T}{\partial a_{ij}}$ is an $n \times n$ matrix in which all entries are zero except the one at j th row and i th column. Thus, we have

$$\text{Trace} \left\{ \left(\mathbf{A}^{-T} \frac{\partial J}{\partial \mathbf{W}} \mathbf{A}^{-T} \right) \frac{\partial \mathbf{A}^T}{\partial a_{ij}} \right\} = \left(\mathbf{A}^{-T} \frac{\partial J}{\partial \mathbf{W}} \mathbf{A}^{-T} \right)_{ij}. \quad (\text{A4})$$

From (A3) and (A4), we can derive

$$\frac{\partial J}{\partial a_{ij}} = \left(\frac{\partial J}{\partial \mathbf{A}} \right)_{ij} = - \left(\mathbf{A}^{-T} \frac{\partial J}{\partial \mathbf{W}} \mathbf{A}^{-T} \right)_{ij}.$$

Hence, we obtain the conclusion

$$\frac{\partial J}{\partial \mathbf{A}} = -\mathbf{A}^{-T} \frac{\partial J}{\partial \mathbf{W}} \mathbf{A}^{-T}.$$

Lemma 2. Let $L(\cdot)$ be a differentiable function with respect to $\mathbf{A} \in \mathbb{R}^{m \times n}$. Consider the following iterative formula

$$\mathbf{A}^{(k+1)} = \mathbf{A}^{(k)} + \mu \nabla L(\mathbf{A}^{(k)}), \quad (\text{A5})$$

where $\nabla L(\mathbf{A}) = \mathbf{U} \mathbf{U}^T \frac{\partial L}{\partial \mathbf{A}}$, and \mathbf{U} is an arbitrary matrix with m rows. If the learning step $\mu > 0$ is sufficiently small, the following inequality holds

$$L(\mathbf{A}^{(k+1)}) \geq L(\mathbf{A}^{(k)}).$$

Proof. Taking the first Taylor extension of $L(\mathbf{A})$, one has

$$L(\mathbf{A} + \boldsymbol{\varepsilon}) = L(\mathbf{A}) + \left\langle \frac{\partial L}{\partial \mathbf{A}}, \boldsymbol{\varepsilon} \right\rangle + o(\boldsymbol{\varepsilon}) = L(\mathbf{A}) + \text{Trace} \left(\left(\frac{\partial L}{\partial \mathbf{A}} \right)^T \cdot \boldsymbol{\varepsilon} \right) + o(\boldsymbol{\varepsilon}). \quad (\text{A6})$$

Substituting $\boldsymbol{\varepsilon} = \mu \nabla L(\mathbf{A}) = \mu \mathbf{U} \mathbf{U}^T \frac{\partial L}{\partial \mathbf{A}}$ into (A6), we obtain

$$\begin{aligned} L(\mathbf{A} + \boldsymbol{\varepsilon}) &= L(\mathbf{A}) + \text{Trace} \left(\left(\frac{\partial L}{\partial \mathbf{A}} \right)^T \cdot \mu \mathbf{U} \mathbf{U}^T \frac{\partial L}{\partial \mathbf{A}} \right) + o(\boldsymbol{\varepsilon}) \\ &= L(\mathbf{A}) + \mu \cdot \text{Trace} \left(\left(\frac{\partial L}{\partial \mathbf{A}} \right)^T \cdot \mathbf{U} \mathbf{U}^T \frac{\partial L}{\partial \mathbf{A}} \right) + o(\boldsymbol{\varepsilon}) \\ &= L(\mathbf{A}) + \mu \cdot \text{Trace} \left(\left(\mathbf{U}^T \frac{\partial L}{\partial \mathbf{A}} \right)^T \cdot \left(\mathbf{U}^T \frac{\partial L}{\partial \mathbf{A}} \right) \right) + o(\boldsymbol{\varepsilon}) \\ &= L(\mathbf{A}) + \mu \cdot \left\| \mathbf{U}^T \frac{\partial L}{\partial \mathbf{A}} \right\|_{\text{Fro}}^2 + o(\boldsymbol{\varepsilon}), \end{aligned} \quad (\text{A7})$$

where $\|\cdot\|_{\text{Fro}}$ is the *Frobenius norm*, and $\text{Trace}(\cdot)$ is the *trace* of a matrix. From (A7), if $\mu > 0$ is sufficiently small, we have $L(\mathbf{A} + \boldsymbol{\varepsilon}) \geq L(\mathbf{A})$ because $\left\| \mathbf{U}^T \frac{\partial L}{\partial \mathbf{A}} \right\|_{\text{Fro}}^2 \geq 0$. Thus, the conclusion

$L(\mathbf{A}^{(k+1)}) = L(\mathbf{A}^{(k)} + \mu \nabla L(\mathbf{A}^{(k)})) \geq L(\mathbf{A}^{(k)})$ is proven.

Remark. Especially, if $\mathbf{U} \neq \mathbf{0}$, $\frac{\partial L}{\partial \mathbf{A}} \neq \mathbf{0}$, we have $L(\mathbf{A}^{(k+1)}) > L(\mathbf{A}^{(k)})$.

- 1 Cao X R, Liu R W. General approach to blind source separation. IEEE Trans Sig Proc, 1996, 44: 562—571
- 2 Comon P. Independent component analysis, a new concept? Sig Proc, 1992, 36(3): 287—314
- 3 Suzuki K, Kiryu T, Nakada T. Fast and precise independent component analysis for high field fMRI time series tailored using prior information on spatiotemporal structure. Human Brain Map, 2001, 15: 54—66
- 4 Lu W, Rajapakse J C. Constrained ICA. NIPS, 2000, 570—576
- 5 Lu W, Rajapakse J C. ICA with reference. Proc ICA2001, 2001, 120—125

- 6 Bell A J, Sejnowski T J. An information-maximization approach to blind source separation and blind deconvolution. *Neural Comput*, 1995, 7: 1129—1159
- 7 Xie S L, Zhang J L. Blind separation algorithm of minimal mutual information based on rotating transform. *Acta Elect Sin* (in Chinese), 2002, 30(5): 628—631
- 8 Hyvarinen A, Oja E. A fast fixed-point algorithm for independent component analysis. *Neural Comput*, 1997, 9(7): 1483—1492
- 9 Tan H Z, Chow T W S. Blind and total identification of ARMA models in higher order cumulants domain. *IEEE Trans Indust Elect*, 1999, 46(6): 1233—1240
- 10 Tan H Z, Chow T W S. Blind identification of quadratic nonlinear models using neural networks with higher order cumulants. *IEEE Trans Indust Elect*, 2000, 47(3): 687—696
- 11 Tan H Z, Aboulnasr T. TOM-based blind identification of nonlinear Volterra systems. *IEEE Trans Instr Meas*, 2006, 55(1): 300—310
- 12 Hyvarinen A. Blind source separation by nonstationarity of variance: a cumulant-based approach. *IEEE Trans Neural Network*, 2001, 12(6): 1471—1474
- 13 Matsuoka K, Ohya M, Kawamoto M. A neural net for blind separation of nonstationary signals. *Neural Networks*, 1995, 8(3): 411—419
- 14 Zhang J L, He Z S, Xie S L. Geometric blind separation algorithm for many source signals. *Chinese J Comput* (in Chinese), 2005, 28(9): 1575—1581
- 15 Li Y Q, Wang J. Blind extraction of singularly mixed source signals. *IEEE Trans Sign Proc*, 2000, 11: 1413—1422
- 16 Li Y Q, Wang J. Sequential blind extraction of instantaneously mixed sources. *IEEE Trans Sig Proc*, 2002, 50(5): 997—1006
- 17 Zibulevsky M, Pearlmutter B A. Blind source separation by sparse decomposition in a signal dictionary. *Neural Comput*, 2001, 13(4): 863—882
- 18 Belouchrani A, Cardoso J F. Maximum likelihood source separation for discrete sources. *Proc EUSIPCO*, 1994, 768—771
- 19 Lee T W, Lewicki M S, Girolami M, et al. Blind source separation of more sources than mixtures using overcomplete representation. *IEEE Sig Proc Lett*, 1999, 6(4): 87—90
- 20 Lewicki M S, Sejnowski T J. Learning overcomplete representations. *Neural Comput*, 2000, 12: 337—365
- 21 Li Y Q, Cichocki A, Amari S. Analysis of sparse representation and blind source separation. *Neural Comput*, 2004, 16(6): 1193—1234
- 22 Bofill P, Zibulevsky M. Underdetermined source separation using sparse representations. *Sig Proc*, 2001, 81: 2353—2362
- 23 Lee T W, Lewicki M, Sejnowski T. ICA mixture models for unsupervised classification of nongaussian sources and automatic context switching in blind signal separation. *IEEE Trans Patt Recogn Mach Intel*, 2000, 22(10): 1—12
- 24 Zhang L Q, Amari S, Cichocki A. Natural gradient approach to blind separation of over- and undercomplete mixtures. In: *Proc of ICA1999*, Aussois, France, Jan. 1999. 455—460
- 25 Donoho D L, Elad M. Maximal sparsity representation via l^1 minimization. In: *Proc National Academy Science*, 2003, 100: 2197—2202
- 26 Takigawa I, Kudo M, Toyama J. Performance analysis of minimum l_1 -norm solutions for underdetermined source separation. *IEEE Trans Sig Proc*, 2004, 52(3): 582—591
- 27 Parra L, Spence C. Convolutional blind separation of nonstationary sources. *IEEE Trans Speech Audio Proc*, 2000, 8: 320—327