

# UNSUPERVISED LEARNING FOR SOURCE SEPARATION WITH MIXTURE OF GAUSSIANS PRIOR FOR SOURCES AND GAUSSIAN PRIOR FOR MIXTURE COEFFICIENTS

Hichem SNOUSSI and Ali MOHAMMAD-DJAFARI  
Laboratoire des Signaux et Systèmes (CNRS – SUPÉLEC – UPS).  
SUPÉLEC, Plateau de Moulon, 91192 Gif-sur-Yvette Cedex, France.  
E-mail: snoussi@lss.supelec.fr, djafari@lss.supelec.fr

**Abstract.** In this contribution, we present two new algorithms for unsupervised learning and source separation for the case of noisy instantaneous linear mixture, within the Bayesian inference framework. The source distribution prior is modeled by a mixture of Gaussians [10] and the mixing matrix elements distributions by a Gaussian. We model the mixture of Gaussians hierarchically by mean of hidden variables representing the labels of the mixture. Then, we consider the joint a posteriori distribution of sources, mixing matrix elements, labels of the mixture and other parameters of the mixture with appropriate prior probability laws to eliminate degeneracy of the likelihood function of variance parameters and we propose two algorithms to estimate sources, mixing matrix and hyperparameters: Joint MAP (Maximum *a posteriori*) algorithm and penalized EM-type algorithm. The performances of these two algorithms are compared through an illustrative example taken in [8].

## PROBLEM DESCRIPTION

We consider the following linear instantaneous mixture of  $n$  sources:

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

where  $\mathbf{x}(t)$  is the  $(m \times 1)$  measurement vector,  $\mathbf{s}(t)$  is the  $(n \times 1)$  source vector which components have to be separated,  $\mathbf{A}$  is the mixing matrix of dimension  $(m \times n)$  and  $\epsilon(t)$  represents noise affecting the measurements. We assume that the  $(m \times T)$  noise matrix  $\epsilon(t)$  is statistically independent of sources, centered, white and Gaussian with covariance matrix  $\mathbf{R}_\epsilon$ . We note  $\mathbf{s}_{1..T}$  the matrix  $n \times T$  of sources and  $\mathbf{x}_{1..T}$  the matrix  $m \times T$  of data.

Source separation problem consists of two sub-problems: Sources restoration and mixing matrix identification. Therefore, three directions can be followed:

1. *Supervised learning*: Identify  $\mathbf{A}$  knowing a training sequence of sources  $\mathbf{s}$ , then use it to reconstruct the sources.
2. *Unsupervised learning*: Identify  $\mathbf{A}$  directly from a part or the whole observations and then use it to recover  $\mathbf{s}$ .
3. *Unsupervised joint estimation*: Estimate jointly  $\mathbf{s}$  and  $\mathbf{A}$

In the following, we investigate the second and third directions. This choice is motivated by practical cases where sources and mixing matrix are unknown.

This paper is organised as follows: We begin in section II by proposing a Bayesian approach to source separation. We set up the notations, present the prior laws of the sources and the mixing matrix elements. We introduce, in section III, a hierarchical modelisation of the sources by mean of hidden variables representing the labels of the mixture of Gaussians in the prior modeling and present the hierarchical JMAP algorithm including estimation of hyperparameters. Since EM algorithm [6] has been used extensively in source separation [3], [1], [2], we considered this algorithm and propose, in section V, a penalized version of the EM algorithm for source separation. This penalization of the likelihood function is necessary to eliminate its degeneracy when some variances of Gaussian mixture approach zero [14], [13], [11]. We will modify the EM algorithm by introducing a classification step and a relaxation strategy to reduce the computational cost. Simulation results are presented in section VI to test and compare the two algorithms performances.

## BAYESIAN APPROACH TO SOURCE SEPARATION

Given the observations  $\mathbf{x}_{1..T}$ , the joint *a posteriori* distribution of unknown variables  $\mathbf{s}_{1..T}$  and  $\mathbf{A}$  is:

$$p(\mathbf{A}, \mathbf{s}_{1..T}, \boldsymbol{\theta} | \mathbf{x}_{1..T}) \propto p(\mathbf{x}_{1..T} | \mathbf{A}, \mathbf{s}_{1..T}, \boldsymbol{\theta}_1) p(\mathbf{A} | \boldsymbol{\theta}_2) p(\mathbf{s}_{1..T} | \boldsymbol{\theta}_3) p(\boldsymbol{\theta}) \quad (2)$$

where  $p(\mathbf{A} | \boldsymbol{\theta}_2)$  and  $p(\mathbf{s}_{1..T} | \boldsymbol{\theta}_3)$  are the prior distributions through which we model our *a priori* information about mixing matrix  $\mathbf{A}$  and sources  $\mathbf{s}$ .  $p(\mathbf{x}_{1..T} | \mathbf{A}, \mathbf{s}_{1..T}, \boldsymbol{\theta}_1)$  is the joint likelihood distribution.  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3)$  are the hyperparameters. From here, we have two directions for unsupervised learning and separation:

1. First, estimate jointly  $\mathbf{s}_{1..T}$ ,  $\mathbf{A}$  and  $\boldsymbol{\theta}$ :

$$(\hat{\mathbf{A}}, \hat{\mathbf{s}}_{1..T}, \hat{\boldsymbol{\theta}}) = \arg \max_{(\mathbf{A}, \mathbf{s}_{1..T}, \boldsymbol{\theta})} \{J(\mathbf{A}, \mathbf{s}_{1..T}, \boldsymbol{\theta}) = \ln p(\mathbf{A}, \mathbf{s}_{1..T}, \boldsymbol{\theta} | \mathbf{x}_{1..T})\} \quad (3)$$

2. Second, integrate (2) with respect to  $\mathbf{s}_{1..T}$  to obtain the marginal in  $(\mathbf{A}, \boldsymbol{\theta})$  and estimate them by:

$$(\hat{\mathbf{A}}, \hat{\boldsymbol{\theta}}) = \arg \max_{(\mathbf{A}, \boldsymbol{\theta})} \{J(\mathbf{A}, \boldsymbol{\theta}) = \ln p(\mathbf{A}, \boldsymbol{\theta} | \mathbf{x}_{1..T})\} \quad (4)$$

Then estimate  $\hat{\mathbf{s}}_{1..T}$  using the posterior  $p(\mathbf{s}_{1..T} | \mathbf{x}_{1..T}, \hat{\mathbf{A}}, \hat{\boldsymbol{\theta}})$ .

The two algorithms we propose follow these two shemes.

### Choice of *a priori* distributions

**Noise *a priori*:** We consider a Gaussian white noise with zero mean and covariance matrix  $\mathbf{R}_\epsilon$  ( $\boldsymbol{\theta}_1 = \mathbf{R}_\epsilon$ ).

**Sources *a priori*:** For sources  $\mathbf{s}$ , we choose a mixture of Gaussians [10]:

$$p(s_j) = \sum_{i=1}^{q_j} \alpha_{ji} \mathcal{N}(m_{ji}, \sigma_{ji}^2), \quad j = 1..n \quad (5)$$

Hyperparameters  $q_j$  are supposed to be known.

This leads to the introduction of hierarchical modelisation  $p(\mathbf{s}_j | z_j) = \mathcal{N}(m_{ji}, \sigma_{ji}^2)$  by considering the hidden variable  $z_j$  taking its values in the discrete set  $\mathcal{Z}_j = (1, \dots, q_j)$  with  $\alpha_{ji} = p(z_j = i)$ .  $\boldsymbol{\theta}_3 = (\alpha_{ji}, m_{ji}, \sigma_{ji}^2)_{j=1..n, i=1..q_j}$ .

**Mixing matrix *a priori*:** To account for some model uncertainty, we assign a Gaussian prior law to each element of the mixing matrix  $\mathbf{A}$ :

$$p(\mathbf{A}_{ij}) = \mathcal{N}(\mathbf{M}_{ji}, \sigma_{a,ij}^2) \quad (6)$$

which can be interpreted as knowing every element ( $\mathbf{M}_{ji}$ ) with some uncertainty ( $\sigma_{a,ij}^2$ ). We underline here the advantage of estimating the mixing matrix  $\mathbf{A}$  and not a separating matrix  $\mathbf{B}$  (inverse of  $\mathbf{A}$ ) which is the case of almost all the existing methods for source separation (see for example [5]). This approach has at least two advantages: (i)  $\mathbf{A}$  does not need to be invertible ( $n \neq m$ ), (ii) naturally, we have some *a priori* information on the mixing matrix not on its inverse which may not exist.

**Hyperparameters *a priori*:** We propose to assign an inverted Gamma prior  $\mathcal{IG}(a, b)$  ( $a > 0$  et  $b > 1$ ) to mixture variances. This prior is necessary to avoid the posterior distribution degeneracy when some variances  $\sigma_{ij}^2$  approche to zero together with noise variance. A more complete study of degeneracies in source separation problem is presented in [14].

## HIERARCHICAL JMAP ALGORITHM

The *a posteriori* distribution of  $\mathbf{s}$  is a mixture of  $\prod_{j=1}^n q_j$  Gaussians. This leads to a high computational cost. To obtain a more reasonable algorithm,

we propose an iterative scalar algorithm by introducing a relaxation procedure: Knowing  $\mathbf{s}_{l \neq j}$ , the *a posteriori* distribution of  $\mathbf{s}_j$  is a mixture of  $q_j$  Gaussians. Including the estimation of hyperparameters, the proposed hierarchical JMAP algorithm follows the following steps in each iteration:

1. Estimate hidden variables  $(\hat{z}_j)_{1..T}$  by:

$$(\hat{z}_j)_{1..T} = (\arg \max_{z_j} p(z_j | \mathbf{x}(t), \hat{\mathbf{A}}, \hat{\mathbf{s}}_{l \neq j}, \hat{\boldsymbol{\theta}}))_{1..T} \quad (7)$$

which permits to estimate partitions:

$$\hat{\mathcal{T}}_{jz} = \{t | (\hat{z}_j)(t) = z\} \quad (8)$$

This corresponds to the classification step.

2. Given the estimate of partitions, hyperparameters  $\hat{m}_{jz}$  and  $\hat{\sigma}_{jz}^2$  are means and variances of Gaussian distributions so the expressions of their posterior estimates are easily derived [15]. Variances are supposed to follow an inverted Gamma prior  $\mathcal{IG}(a, b)$ . The hyperparameter  $\hat{\alpha}_{jz}$  is updated as:

$$\hat{\alpha}_{jz} = \text{Card}(\hat{\mathcal{T}}_{jz})/T \quad (9)$$

3. Estimation of sources using  $\hat{\mathbf{s}}_{1..T} = \arg \max_{\mathbf{s}_{1..T}} \{p(\mathbf{s}_{1..T} | \mathbf{x}_{1..T}, \hat{\mathbf{A}}, \hat{\boldsymbol{\theta}})\}$ .
4. Estimation of mixing matrix using  $\hat{\mathbf{A}} = \arg \max_{\mathbf{A}} \{p(\mathbf{A} | \mathbf{x}_{1..T}, \hat{\mathbf{s}}_{1..T}, \hat{\boldsymbol{\theta}})\}$ .

### Penalized EM-type Algorithm

The EM algorithm has been used extensively in data analysis to find the maximum likelihood estimation of a set of parameters from given data [12], [6], [7]. Considering both the mixing matrix  $\mathbf{A}$  and hyperparameters  $\boldsymbol{\theta}$ , at the same level, being unknown parameters and complete data being  $\mathbf{x}_{1..T}$  and  $\mathbf{s}_{1..T}$ , the EM algorithm writes: (i) E-step (expectation) consists in forming the logarithm of the joint distribution of observed data  $\mathbf{x}$  and hidden data  $\mathbf{s}$  conditionally to parameters  $\mathbf{A}$  and  $\boldsymbol{\theta}$  and then compute its expectation conditionally to  $\mathbf{x}$  and estimated parameters  $\mathbf{A}'$  and  $\boldsymbol{\theta}'$  (evaluated in the previous iteration), (ii) M-step (maximization) consists of the maximization of the obtained functional with respect to the parameters  $\mathbf{A}$  and  $\boldsymbol{\theta}$ .

Recently, in [3], [1] an EM algorithm has been used in source separation with mixture of Gaussians as sources prior. In this work, we show that:

1. This algorithm fails in estimating jointly variances of Gaussian mixture and noise covariance matrix. We proved that this is due to the degeneracy of the estimated variance to zero and a problem of identifiability.
2. The computational cost of this algorithm is very high.

3. The algorithm is very sensitive to initial conditions.
4. In [3], there's neither an *a priori* distribution on the mixing matrix  $\mathbf{A}$  nor on the hyperparameters  $\boldsymbol{\theta}$ .

Here, we propose to extend this algorithm in two ways by:

1. Introducing an *a priori* distribution for  $\boldsymbol{\theta}$  to eliminate degeneracy. This *a priori* contributes in reducing the problem of non identifiability but doesn't eliminate it completely.
2. Introducing an *a priori* distribution for  $\mathbf{A}$  to express our previous knowledge on the mixing matrix.
3. Taking advantage of our hierarchical model and the idea of classification to reduce the computational cost.

To distinguish the proposed algorithm from the one proposed in [3], we call this algorithm the *Penalized EM* algorithm. The two steps then become:

1. E-step :  $Q(\mathbf{A}, \boldsymbol{\theta} | \mathbf{A}', \boldsymbol{\theta}') = E_{\mathbf{x}, \mathbf{s}}[\log p(\mathbf{x}, \mathbf{s} | \mathbf{A}, \boldsymbol{\theta}_1, \boldsymbol{\theta}_3) + \log p(\mathbf{A} | \boldsymbol{\theta}_2) + \log p(\boldsymbol{\theta} | \mathbf{x}, \mathbf{A}', \boldsymbol{\theta}')] ]$
2. M-step :  $(\hat{\mathbf{A}}, \hat{\boldsymbol{\theta}}) = \arg \max_{(\mathbf{A}, \boldsymbol{\theta})} Q(\mathbf{A}, \boldsymbol{\theta} | \mathbf{A}', \boldsymbol{\theta}')$

We suppose in the following that  $(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$  are known (noise variance and mixing matrix *a priori* parameters). The joint distribution is factorized as:  $p(\mathbf{x}, \mathbf{s}, \mathbf{A}, \boldsymbol{\theta}) = p(\mathbf{x} | \mathbf{A}, \mathbf{s}, \boldsymbol{\theta}_1) p(\mathbf{A} | \boldsymbol{\theta}_2) p(\mathbf{s} | \boldsymbol{\theta}_3) p(\boldsymbol{\theta}_3)$ . We can remark that  $p(\mathbf{x}, \mathbf{s}, \mathbf{A}, \boldsymbol{\theta})$  as a function of  $(\mathbf{A}, \boldsymbol{\theta}_3)$  is separable in  $\mathbf{A}$  and  $\boldsymbol{\theta}_3$ . Consequently, the functional is separated into two factors: one representing an  $\mathbf{A}$  functional and the other representing a  $\boldsymbol{\theta}_3$  functional:

$$Q(\mathbf{A}, \boldsymbol{\theta}_3 | \mathbf{A}', \boldsymbol{\theta}'_3) = Q_a(\mathbf{A} | \mathbf{A}', \boldsymbol{\theta}'_3) + Q_h(\boldsymbol{\theta}_3 | \mathbf{A}', \boldsymbol{\theta}'_3) \quad (10)$$

- **Maximization with respect to  $\mathbf{A}$ :** By introducing the Kronecker product [4], we can derive an explicit expression of the update of  $\mathbf{A}$  maximizing the  $Q_a$  functional:

$$\mathbf{Vec}(\mathbf{A}) = \left[ T \hat{\mathbf{R}}'_{ss} \otimes \mathbf{R}_\epsilon^{-1} + \text{diag}(\text{Vec}(\boldsymbol{\Gamma})) \right]^{-1} \text{Vec}(T \mathbf{R}_\epsilon^{-1} \hat{\mathbf{R}}_{xs} + \boldsymbol{\Gamma} \odot \mathbf{M}) \quad (11)$$

where  $\otimes$  is the Kronecker product,  $\odot$  is the element-by-element product and  $\mathbf{Vec}(\cdot)$  is the column presentation of a matrix.  $\boldsymbol{\Gamma}$  is the matrix  $(1/\sigma_{a,ij}^2)$  and  $(\hat{\mathbf{R}}_{xs}, \hat{\mathbf{R}}_{ss})$  are the following statistics:

$$\begin{cases} \hat{\mathbf{R}}_{xs} &= \frac{1}{T} \sum_{t=1}^T E[\mathbf{x}(t) \mathbf{s}(t)^T | \mathbf{x}, \mathbf{A}', \boldsymbol{\theta}'] \\ \hat{\mathbf{R}}_{ss} &= \frac{1}{T} \sum_{t=1}^T E[\mathbf{s}(t) \mathbf{s}(t)^T | \mathbf{x}, \mathbf{A}', \boldsymbol{\theta}'] \end{cases} \quad (12)$$

Evaluation of  $\hat{\mathbf{R}}_{xs}$  and  $\hat{\mathbf{R}}_{ss}$  requires the computation of the expectations of  $\mathbf{x}(t) \mathbf{s}(t)^T$  and  $\mathbf{s}(t) \mathbf{s}(t)^T$ . The main computational cost is due to the fact

that the expectation of any function  $f(\mathbf{s})$  is given by:

$$E[f(\mathbf{s}) | \mathbf{x}, \mathbf{A}', \boldsymbol{\theta}'] = \sum_{\mathbf{z}' \in \prod_{i=1}^n \mathcal{Z}_i} E[f(\mathbf{s}) | \mathbf{x}, \mathbf{z} = \mathbf{z}', \mathbf{A}', \boldsymbol{\theta}'] p(\mathbf{z}' | \mathbf{x}, \mathbf{A}', \boldsymbol{\theta}'). \quad (13)$$

which involves a sum of  $\prod_{j=1}^n q(j)$  terms corresponding to the whole combinations of labels. One way to obtain an approximate but fast estimate of this expression is to limit the summation to only one term corresponding to the MAP estimate of  $\mathbf{z}$ :

$$E[f(\mathbf{s}) | \mathbf{x}, \mathbf{A}', \boldsymbol{\theta}'] = E[f(\mathbf{s}) | \mathbf{x}, \mathbf{z} = \hat{\mathbf{z}}^{MAP}, \mathbf{A}', \boldsymbol{\theta}'] .$$

**Maximisation with respect to  $\boldsymbol{\theta}_3$ :** With an uniform *a priori* for the means and variances, maximisation of the functional  $Q$  with respect to  $\boldsymbol{\theta}_3$  gives :

$$\begin{aligned} \hat{\alpha}_{jz} &= \frac{\sum_{t=1}^T p(z_j(t) | \mathbf{x}, \mathbf{A}', \boldsymbol{\theta}')}{T} \\ \hat{m}_{jz} &= \frac{\sum_{t=1}^T \mu_{jz}(t) p(z_j(t) | \mathbf{x}, \mathbf{A}', \boldsymbol{\theta}')}{\sum_{t=1}^T p(z_j(t) | \mathbf{x}, \mathbf{A}', \boldsymbol{\theta}')} \\ \hat{\sigma}_{jz}^2 &= \frac{\sum_{t=1}^T (V_{jz}(t) + \mu_{jz}^2(t) - 2\hat{m}_{jz}\mu_{jz}(t) + \hat{m}_{jz}^2) p(z_j(t) | \mathbf{x}, \mathbf{A}', \boldsymbol{\theta}')}{\sum_{t=1}^T p(z_j(t) | \mathbf{x}, \mathbf{A}', \boldsymbol{\theta}')} \end{aligned}$$

where:

$$\begin{aligned} \mu_{jz}(t) &= \mathbf{E}[s_j(t) | \mathbf{x}(t), z] \\ V_{jz}(t) &= \mathbf{E}[s_j(t)^2 | \mathbf{x}(t), z] \end{aligned}$$

The computation of  $p(z_j(t) | \mathbf{x}, \mathbf{A}', \boldsymbol{\theta}')$  needs a summation over all combinations of labels:

$$p(z_j(t) | \mathbf{x}, \mathbf{A}', \boldsymbol{\theta}') = \sum_{\mathbf{z} \in \mathcal{Z} | z(j) = z_j(t)} p(\mathbf{z} | \mathbf{x}(t), \mathbf{A}', \boldsymbol{\theta}') \quad (14)$$

The relaxation strategy consists on replacing expression (14) by:

$$p(z_j(t) | \mathbf{x}, \mathbf{A}', \boldsymbol{\theta}', \hat{s}_{l \neq j})$$

which is obtained by integrating only with respect to  $s_j$ , the other components are fixed and set to their MAP estimates in the previous iteration. Assigning an Inverted Gamma prior  $\mathcal{IG}(a, b)$  ( $a > 0$  et  $b > 1$ ) to the variances, the re-estimation equations become:

$$\hat{\alpha}_{jz} = \frac{\sum_{t=1}^T p(z_j(t) | \mathbf{x}, \hat{s}_{l \neq j})}{T} \quad (15)$$

$$\hat{m}_{jz} = \frac{\sum_{t=1}^T \mu_{jz}(t) p(z_j(t) | \mathbf{x}(t), \hat{s}_{l \neq j})}{\sum_{t=1}^T p(z_j(t) | \mathbf{x}, \hat{s}_{l \neq j})} \quad (16)$$

$$\hat{\sigma}_{jz}^2 = \frac{2b + \sum_{t=1}^T (V_{jz} + \mu_{jz}^2 - 2\hat{m}_{jz}\mu_{jz} + \hat{m}_{jz}^2) p(z(t) | \mathbf{x}, \hat{s}_{l \neq j})}{\sum_{t=1}^T p(z(t) | \mathbf{x}, \hat{s}_{l \neq j}) + 2(a-1)} \quad (17)$$

**Summary of the penalized EM-type-type algorithm** Based on the preceding equations, we propose the following algorithm to estimate sources and parameters using the following five steps:

1. Update of data classification by estimating  $\hat{z}_{1..T}$  using 7 as in JMAP.
2. Estimate the mixing matrix  $\mathbf{A}$  according to the re-estimation equation (11).
3. Given this classification, sources estimate is the mean of the Gaussian *a posteriori* law.
4. Estimate the hyperparameters according to (15), (16) and (17).

## SIMULATION RESULTS

To be able to compare the results obtained by the two proposed algorithms with the results obtained by some other classical methods, we generated data according to the example described in [8].

**Data generation:** 2 sources, each component *a priori* is a mixture of two Gaussians ( $\pm 1$ ),  $\psi = 1/\sigma^2 = 100$  for all Gaussians. These original sources are mixed with the mixing matrix  $\mathbf{A} = \begin{pmatrix} 1 & -0.6 \\ 0.4 & 1 \end{pmatrix}$ . A noise of variance  $\sigma_\epsilon^2 = 0.03$  is added ( $SNR = 15$  dB). Number of observations is 1000.

**Parameters:**  $\mathbf{M} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ ,  $\mathbf{\Gamma} = (1/\sigma_{a,ij}^2) = \begin{pmatrix} 150 & 0.009 \\ 0.009 & 150 \end{pmatrix}$ ,

$\mathbf{\Pi} = (\alpha_{jz}) = \begin{pmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{pmatrix}$ ,  $a = 200$  and  $b = 2$ .

**Initial conditions:**  $\mathbf{A}^{(0)} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ ,  $\psi^{(0)} = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$ ,  $m^{(0)} = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$

and  $\mathbf{s}^{(0)}$  generated according to  $\mathbf{s}_j^{(0)} \sim \sum_{z=1}^{q_j} \Pi_{jz} \mathcal{N}(m_{jz}^{(0)}, \frac{1}{\psi_{jz}^{(0)}})$ .

**Results with JMAP algorithm:.** Sources are recovered with negligible mean quadratic error:  $MEQ(\mathbf{s}_1) = 0.0094$  and  $MEQ(\mathbf{s}_2) = 0.0097$ . The

non-negative performance index of [9] is used to characterize mixing matrix identification achievement:

$$ind(S = \hat{\mathbf{A}}^{-1} \mathbf{A}) = \frac{1}{2} \left[ \sum_i \left( \sum_j \frac{|S_{ij}|^2}{\max_l |S_{il}|^2} - 1 \right) + \sum_j \left( \sum_i \frac{|S_{ij}|^2}{\max_l |S_{lj}|^2} - 1 \right) \right]$$

Figure 1a represents the index evolution through iterations. Note the convergence of JMAP algorithm since iteration 30 to a satisfactory value of  $-45 \text{ dB}$ . For the same SNR, algorithms PWS, NS [8] and EASI [5] reach a value greater than  $-35 \text{ dB}$  after 6000 observations. Figures 1b and 1c illustrate the identification of hyperparameters. We note the convergence of the parameters to the original values ( $-1$  for  $m_{11}$  and  $100$  for  $\psi_{11}$ ). In order to validate the idea of data classification before estimating hyperparameters, we can visualize the evolution of classification error (number of data badly classified). Figure 1d shows that this error converges to zero at iteration 15. Then, after this iteration, hyperparameters identification is performed with the right classified data: estimation of  $m_{jz}$  and  $\psi_{jz}$  uses only data which belong to this class and is not corrupted by other data which bring erroneous information on these hyperparameters.

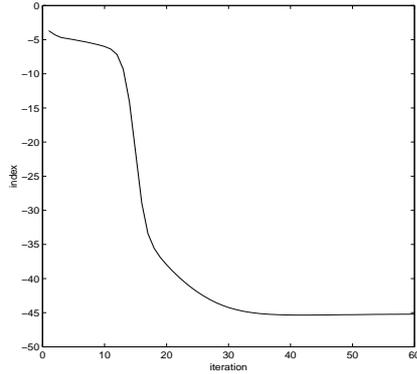


Figure 1-a- Evolution of index

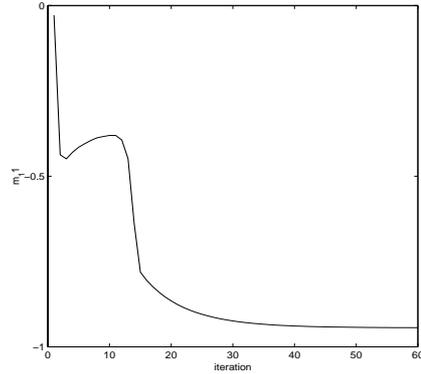


Figure 1-b- Identification of  $m_{11}$

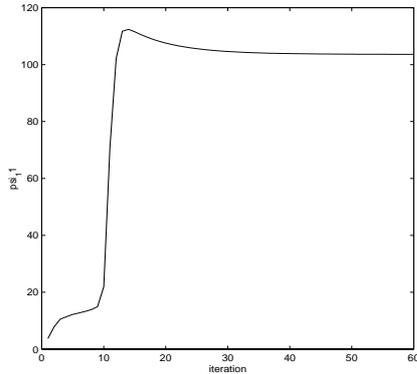


Figure 1-c- Identification of  $\psi_{11}$

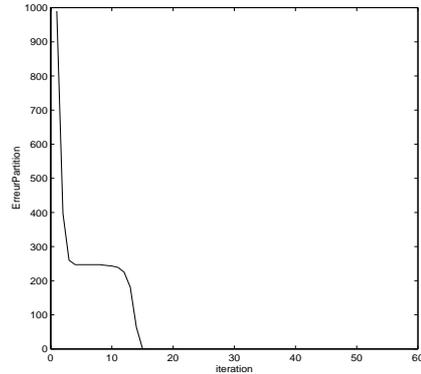


Figure 1-d- Evolution of classification error

**Results with Penalized EM-type algorithm:** The penalized EM-type algorithm has an optimization cost approximately 2 times higher, per sample, than the JMAP algorithm. However, both algorithms have a reasonable computational complexity, linearly increasing with the number of samples. Sensitivity to initial conditions is inherent to the EM-algorithm even to the penalized version. In order to illustrate this fact, we simulated the algorithm with the same parameters as above. Note that initial conditions for hyperparameters are  $\psi^{(0)} = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$  and  $m^{(0)} = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$ . However, the penalized EM-type algorithm fails in separating sources. We note then that JMAP algorithm is more robust to initial conditions. We modified the initial condition to have means:  $m^{(0)} = \begin{pmatrix} -0.5 & 0.5 \\ -0.5 & 0.5 \end{pmatrix}$ . We noted, in this case, the convergence of the penalized EM-type algorithm to the correct solution. Figures 2-a and 2-b illustrate the separation results:

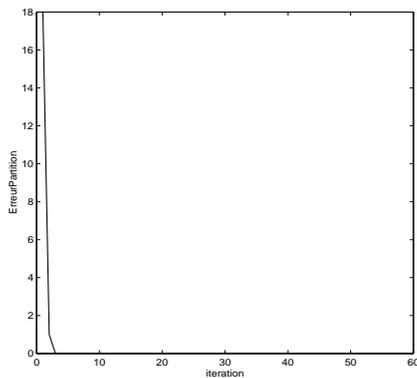


Figure 2-a- Evolution of classification error

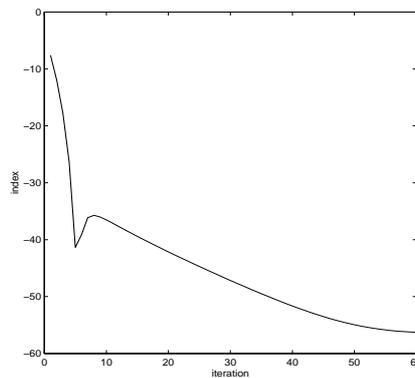


Figure 2-b- Evolution of index

## CONCLUSION

We proposed two new algorithms for unsupervised learning and source separation when the sources distributions are modeled to be a mixture of Gaussians. Considering the mixture model as a hierarchical modeling with hidden variables representing labels, we introduced a classification step before the estimation of hyperparameters. This classification step is useful, not only to do a better job in the estimation of the mixing components parameters, but also to reduce the computational cost of JMAP and Penalized EM algorithms.

It is also important to mention that the Bayesian estimation framework we have adopted has specific aspects including the introduction of a *a priori* distribution for the mixing matrix and hyperparameters. This was motivated by two different reasons: Mixing matrix prior should exploit previous information and variances prior should regularize the log-posterior objective function.

## REFERENCES

- [1] A. Belouchrani, **Séparation autodidacte de sources: Algorithmes, Performances et Application à des signaux expérimentaux**, Phd thesis, Ecole Nationale Supérieure des Télécommunications, 1995.
- [2] A. Belouchrani and J.-F. Cardoso, "Maximum likelihood source separation for discrete sources," in **EUSIPCO'94**, 1994.
- [3] O. Bermond, **Méthodes statistiques pour la séparation de sources**, Phd thesis, Ecole Nationale Supérieure des Télécommunications, 2000.
- [4] J. W. Brewer, "Kronecker products and matrix calculus in system theory," **IEEE Trans. Circ. Syst.**, vol. CS-25, no. 9, pp. 772–781, 1978.
- [5] J. Cardoso and B. Labeld, "Equivariant adaptative source separation," **Signal Processing**, vol. 44, pp. 3017–3030, 1996.
- [6] A. P. Dempster, N. M. Laird and D. B. Rubin, "Maximum Likelihood from incomplete data via the EM algorithm," **J. R. Statist. Soc. B**, vol. 39, pp. 1–38, 1977.
- [7] A. O. Hero and J. A. Fessler, "Asymptotic Convergence Properties of EM-Type Algorithms," Preprints 85-T-21, **Dept. of Electrical Engineering and Computer Science, University of Michigan**, 1985.
- [8] O. Macchi and E. Moreau, "Adaptative unsupervised separation of discrete sources," in **Signal Processing**, 1999, vol. 73, pp. 49–66.
- [9] E. Moreau and O. Macchi, "High-order contrasts for self-adaptative source separation," in **Adaptative Control Signal Process.** 10, 1996, pp. 19–46.
- [10] E. Moulines, J. Cardoso and E. Gassiat, "Maximum likelihood For Blind Separation And Deconvolution Of Noisy Signals Using Mixture Models," in **ICassp-97**, Munich, Germany, April 1997.
- [11] D. Ormoneit and V. Tresp, "Averaging, Maximum Penalized Likelihood and Bayesian Estimation for Improving Gaussian Mixture Probability Density Estimates," **IEEE Transactions on Neural Networks**, vol. 9, no. 4, pp. 639–649, July 1998.
- [12] R. A. Redner and H. F. Walker, "Mixture densities, maximum likelihood and the EM algorithm," **SIAM Rev.**, vol. 26, no. 2, pp. 195–239, April 1984.
- [13] A. Ridolfi and J. Idier, "Penalized Maximum Likelihood estimation for Univariate Normal Mixture Distributions," in **Actes 17<sup>e</sup> coll. GRETSI**, Vannes, France, September 1999, pp. 259–262.
- [14] H. Snoussi and A. Mohammad-Djafari, "Dégénérescences des estimateurs MV en séparation de sources," Technical report ri-s0010, **gpi-12s**.
- [15] H. Snoussi and A. Mohammad-Djafari, "Bayesian source separation with mixture of Gaussians prior for sources and Gaussian prior for mixture coefficients," in A. Mohammad-Djafari (ed.), **Bayesian Inference and Maximum Entropy Methods**, MaxEnt Workshops, Gif-sur-Yvette, France: to appear in Amer. Inst. Physics, July 2000.