

An Expectation-Maximization Algorithm for Blind Separation of Noisy Mixtures Using Gaussian Mixture Model

Fanglin Gu, Hang Zhang, *Member, IEEE*, Wenwu Wang, *Senior Member, IEEE*, and Chunlin Xiong

Abstract—In this paper, we propose a new expectation-maximization (EM) algorithm, named GMM-EM, to blind separation of noisy instantaneous mixtures, in which the non-Gaussianity of independent sources is exploited by modeling their distribution using the Gaussian mixture model (GMM). The compatibility between the incomplete data structure of the GMM and the hidden variable nature of the source separation problem leads to an efficient hierarchical learning and alternative method for estimating the sources and the mixing matrix. In comparison with conventional blind source separation algorithms, the proposed GMM-EM algorithm has superior performance for the separation of noisy mixtures due to the fact that the covariance matrix of the additive Gaussian noise is treated as a parameter. Furthermore, the GMM-EM algorithm works well in underdetermined cases by incorporating any prior information one may have and jointly estimating the mixing matrix and source signals in a Bayesian framework. Systematic simulations with both synthetic and real speech signals are used to show the advantage of the proposed algorithm over conventional independent component analysis techniques, such as FastICA, and a recent technique called null space component analysis (NCA), especially for noisy and/or underdetermined mixtures.

Index Terms—Blind source separation, Gaussian mixture model, expectation-maximization, underdetermined mixture.

I. INTRODUCTION

BLIND source separation (BSS) aims to estimate unknown sources from the observed sensor signals without (or with very limited) prior information about the sources and how the sources propagate to the sensors. It has drawn great attention due to its wide range of applications in signal processing. Many algorithms have been developed to solve the BSS problem based on the assumption that the sources to be recovered are statistically independent, leading to a family of

well-known methods called independent component analysis (ICA) [1]–[5], such as, the information maximization based Infomax algorithm [1], the joint approximate diagonalization of eigenmatrices (JADE) algorithm [3] and the FastICA algorithm [4]. However, most of the ICA methods are developed for the case of determined/overdetermined mixtures (i.e., the number of sources P is equal to or smaller than the number of sensors Q), and consider a noiseless source separation model. These methods are not directly applicable to the problem of source separation from noisy and/or underdetermined mixtures.

There have been a number of attempts to extend the ICA approach in order to address the noisy and/or underdetermined BSS problem, such as, the noise-model based FastICA algorithm [6], cumulant-based separation algorithms [7], [8], and the characteristic functions based blind identification methods [9]–[12]. These methods have offered new ideas for estimating the mixing system, however, their performance is still limited for recovering the source signals. The underdetermined source separation problem is, in particular, very challenging since, as opposed to dealing with determined/over-determined mixtures, even provided with the information about the mixing system, the sources cannot be uniquely reconstructed, simply because, for $P > Q$, the mixing matrix is not invertible. In contrast to the ICA approach, the recently proposed null space component analysis (NCA) approach can solve the noisy and/or underdetermined BSS problem effectively [13], [14]. Given a set of signals, the NCA approach constructs an operator for each signal so that only the signal of interest is in the operator's null space, and all the other signals are excluded. Furthermore, an additional constraint on the rank of the operators is imposed to remove the rotation ambiguity.

In fact, the methods discussed above can be considered as special cases under the Bayesian framework. In a Bayesian technique, a statistical model defined by a set of parameters is used to describe the source separation problem [15]–[18]. The parameters of the model can be inferred from the acquired data, with the help of some prior information about the physical system under consideration. As compared with the classical ICA methods [1]–[5], the Bayesian approach provides advantages in several scenarios. For example, the Bayesian approach is often much more robust to noise since the noise levels in the data are taken into account through the parameterization of the noise covariance matrix within the Bayesian model [15], [16]. Second, the Bayesian approach also enables any prior knowledge about the physical application systems to be exploited in the model where appropriate prior

This work is supported in part by the Natural Science Foundation of China under Grant 61001106, the National Program on Key Basic Research Project of China under Grant 2009CB320400, the Major Projects of the National Natural Science Foundation of China under Grant 91338105, and the foundation of Science and Technology on Information Transmission and Dissemination in Comm. Networks Lab.

F. Gu and C. Xiong are with the School of Electronic Science and Engineering, National University of Defense Technology, Changsha 410073, P. R. China. C. Xiong is also with the Science and Technology on Information Transmission and Dissemination in Communication Networks Laboratory, Shijiazhuang 050081, P. R. China. e-mail: (gu.fanglin@gmail.com, xch-lzju@nudt.edu.cn).

H. Zhang is with the College of Communication Engineering, PLA University of Science and Technology, Nanjing 210007, P. R. China. e-mail: (hangzh_2002@163.com).

W. Wang is with the Department of Electronic Engineering, University of Surrey, Guildford GU2 7XH, U. K. e-mail: (w.wang@surrey.ac.uk)

distributions for the unknown parameters can be assigned. Moreover, the BSS problem can be reformulated as a problem of joint maximum a posteriori (MAP) probability estimation of the mixing matrix and the sources, and as a result, the Bayesian approach can be extended to address the underdetermined BSS problem [17], [18].

Under the Bayesian framework, a number of BSS approaches have been presented in the literature. Belouchrani *et al.* [19] developed a maximum-likelihood (ML) method for jointly estimating the mixing matrix and noise covariance matrix via the expectation-maximization (EM) algorithm [20], [21] where the sources are drawn from a finite alphabet set. However, many natural signals, such as speech signals, are continuous signals (rather than the discrete sources). It has been shown in [22] that many PDFs can be closely approximated by a finite-order Gaussian mixture model (GMM) via the Kullback-Leibler (KL) divergence [23]. Following this route, several GMM based BSS approaches have been proposed. For example, an approximate ML approach was developed by Moulines *et al.* [24] for blind separation and deconvolution of noisy linear mixtures. This method primarily considered the use of GMMs to model the distribution of sources, and the parameters of this model, along with the unknown mixing matrix were optimized to best represent the observed data. Some related works can also be found in [15], [17], [18], [25]–[27], in which different mixture models, such as generalized Gaussian mixture model [25], are used. The application of these models for blind separation of underdetermined mixtures has been studied in [17], [18], [25]. On the other hand, an alternative strategy, in which the GMM is fitted to the observed data, rather than the sources, is developed in [28], [29]. In this approach, the mixtures are separated by finding the rotation matrix that approximately diagonalizes all of the correlation matrices resulting from the GMM. However, it is limited to the noiseless determined case.

Recently, we proposed an EM algorithm for separating noisy determined/underdetermined mixtures with non-stationary sources, in which the continuous density hidden Markov model (CDHMM) is used to model the PDF and to track the non-stationarity of the sources [30]. Preliminary study on synthesized data has shown great potentials of this algorithm, despite the challenge in initialising appropriately the large number of hyper-parameters in practical scenarios. In practice, the distribution model plays a vital role in Bayesian approach. On one hand, the distribution model is required to depict as many distribution forms and traits as possible to make the approach more flexible and to potentially improve separation performance. On the other hand, it is also required to involve as few parameters as possible such that the Bayesian approach can be implemented easily. Hence, there is a tradeoff between the generalization of the distribution model and estimation precision.

In this paper, we also consider the challenging noisy and/or underdetermined BSS problem. In order to address the above issue, we propose to exploit the non-Gaussianity of the sources by modeling their distributions using a GMM, and to incorporate prior information by assigning conjugate priors for the parameters of the GMM and mixing matrix for improving the

separation performance. In such a case, the BSS problem can be treated as a problem of estimating parameters from incomplete data. The EM algorithm is probably the most well-known algorithm for obtaining the ML estimates in parametric models for incomplete data. It is an iterative algorithm alternating between the E-step and M-step respectively. In the E-step, the conditional expectation of the complete data log-likelihood is computed on the basis of the observed data and parameter estimates. In the M-step, the parameters are estimated by maximizing the complete data log-likelihood from the E-step. Therefore, an EM algorithm is proposed for obtaining the MAP estimates of the mixing matrix, the sources and the noise covariance matrix in a joint manner. Although there are some similar works in the literature, our approach differs from these works in the following aspects. First, such as [24], the conjugate priors used for incorporating prior information have not been considered. Second, as opposed to the variational Bayesian method in [27], a new GMM-EM method is used to obtain the MAP estimates of the sources and parameters due to its advantage of providing fast and stable convergence [31]. Thanks to the prior information incorporated by the conjugate priors, the proposed EM method works well even for underdetermined mixtures. Third, in comparison with the method based on CDHMM [30], the proposed GMM-EM method is easier to implement.

The remainder of this paper is organized as follows. In Section II, the BSS problem and the assumptions made in our work are presented. The source distribution model based on GMM is given in Section III. The notations describing the prior laws for the mixing coefficients, noise covariance matrix and the model parameters are presented in Section IV. In Section V, a new GMM-EM algorithm is derived for the estimation of the mixing coefficients, the noise covariance matrix, and the model parameters, in order to estimate the source signals. Issues regarding the practical implementation and performance of the proposed algorithm are discussed in Section VI, where the initialization scheme for the parameters, the convergence performance, and computational complexity are analyzed. In Section VII, simulations are provided to show the performance of the proposed algorithm. Finally, conclusions are drawn in Section VIII.

II. PROBLEM FORMULATION

We consider the well-known instantaneous linear mixing model given as [4]

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

The random vector $\mathbf{s}(t) = [s_1(t), \dots, s_P(t)]^T$, representing P statistically independent sources at discrete time instance t , is mixed by a time-invariant unknown mixing matrix \mathbf{A} . The observation vector $\mathbf{x}(t) = [x_1(t), \dots, x_Q(t)]^T$ is obtained from an array of Q sensors, and contaminated by the noise vector $\mathbf{w}(t) = [w_1(t), \dots, w_Q(t)]^T$ which is assumed to be Gaussian white with zero-mean and unknown covariance matrix $\mathbf{R}_w = \text{diag}(\sigma_1^2, \dots, \sigma_Q^2)$ and independent of $\mathbf{s}(t)$.

Different from the determined and/or noiseless models investigated in many existing contributions, here, we consider

the more practical situations where the mixtures may be corrupted by noise or the mixing system is underdetermined. Our objective in this paper is therefore to develop an algorithm for recovering the source signals from noisy and/or underdetermined mixtures. To this aim, we propose to reconstruct the source signals $\{\mathbf{s}(t)\}_{t=1,\dots,T}$ and the mixing matrix \mathbf{A} in a joint manner under the Bayesian framework on the basis of the observed signals $\{\mathbf{x}(t)\}_{t=1,\dots,T}$ and the assignment of some prior information. Note that the method described in this paper could be extended to a convolutive mixing model, which however is out of the scope of this work.

III. SOURCE DISTRIBUTION MODEL

This section describes the source model based on GMM. The PDF of the i th source signal at time instance t is modeled by the GMM as follows

$$f_s(s_i(t), \theta_i) = \sum_{l_i=1}^{N_i} \alpha_{i,l_i} \mathcal{N}(s_i(t); \mu_{i,l_i}, \sigma_{i,l_i}^2), i = 1, \dots, P \quad (2)$$

where $\mathcal{N}(\cdot; \cdot, \cdot)$ denotes a Gaussian density function and N_i denotes the number of Gaussians. The mixing weights are denoted by $\{\alpha_{i,l_i}\}_{l_i}^{N_i}$, such that $\sum_{l_i=1}^{N_i} \alpha_{i,l_i} = 1$. The means and variances of the Gaussians are denoted by $\{\mu_{i,l_i}\}_{l_i}^{N_i}$ and $\{\sigma_{i,l_i}^2\}_{l_i}^{N_i}$, respectively. Assuming that the source signals are statistically independent, the joint PDF of the sources can be formulated as follows [26]

$$\begin{aligned} f_s(\mathbf{s}(t); \Theta) &= \prod_{i=1}^P f_s(s_i(t); \theta_i) \\ &= \sum_{l_1}^{N_1} \alpha_{1,l_1} \mathcal{N}(s_1(t); \mu_{1,l_1}, \sigma_{1,l_1}^2) \\ &\quad \sum_{l_2}^{N_2} \alpha_{2,l_2} \mathcal{N}(s_2(t); \mu_{2,l_2}, \sigma_{2,l_2}^2) \\ &\quad \cdots \sum_{l_P}^{N_P} \alpha_{P,l_P} \mathcal{N}(s_P(t); \mu_{P,l_P}, \sigma_{P,l_P}^2) \\ &= \sum_{l_1}^{N_1} \sum_{l_2}^{N_2} \cdots \sum_{l_P}^{N_P} \omega_{l_1, l_2, \dots, l_P} \\ &\quad \mathcal{N}\left([s_1(t), s_2(t) \cdots, s_P(t)]^T; \right. \\ &\quad \left. [\mu_{1,l_1}, \mu_{2,l_2}, \dots, \mu_{P,l_P}]^T, \right. \\ &\quad \left. \text{diag}(\sigma_{1,l_1}^2, \sigma_{2,l_2}^2, \dots, \sigma_{P,l_P}^2)\right) \\ &= \sum_{m=1}^M \omega_m \mathcal{N}(\mathbf{s}(t); \boldsymbol{\mu}_m, \mathbf{C}_m) \end{aligned} \quad (3)$$

where $M = \prod_{i=1}^P N_i$ is the total number of Gaussians in the joint PDF and $\omega_m = \prod_{i=1}^P \alpha_{i,l_i}$; $m = 1, \dots, M$ are the mixing weights of each Gaussian component such that $\sum_{m=1}^M \omega_m = 1$. The index denotes a unique combination of the Gaussian components from each source, i.e., $l_1, \dots, l_P \rightarrow m$, where $l_i \in \{1, \dots, N_i\}$ denotes a Gaussian index of the source. The mean vector and covariance matrix of the Gaussian are denoted by $\boldsymbol{\mu}_m = [\mu_{1,l_1}, \mu_{2,l_2}, \dots, \mu_{P,l_P}]^T$ and

$\mathbf{C}_m = \text{diag}(\sigma_{1,l_1}^2, \sigma_{2,l_2}^2, \dots, \sigma_{P,l_P}^2)$, respectively. It can be observed from (3) that the joint PDF of the sources is a multivariate GMM parameterized by a diagonal covariance matrix [30].

IV. CHOICES OF PRIOR DENSITIES

In this section, we discuss how to choose the prior distributions to incorporate prior information for improving the performance of blind separation. The prior distribution is used to attribute uncertainty rather than randomness to the unknown parameter or latent variable. The hyper-parameters of the priors are chosen to reflect any existing information. In the Bayesian framework, the aim for solving BSS problem is to obtain the posterior distribution of the relevant parameters. Generally, let z be a random variable, and ϑ be the relevant parameter. According to the Bayesian theorem, the posterior distribution can be represented as the product of the likelihood function $f(z|\vartheta)$ and prior $f(\vartheta)$, normalized by the probability of the data $f(z)$

$$f(\vartheta|z) = \frac{f(z|\vartheta)f(\vartheta)}{\int f(z|\vartheta)f(\vartheta)d\vartheta} \quad (4)$$

The likelihood function is usually well-determined from the data-generating process and can be considered fixed. Therefore, the difficulty in calculating the integral in the denominator in the right hand side (RHS) of the above equation will depend on the choice of the prior distributions. With conjugate priors¹, the posterior distribution will have the same algebraic form as the prior distribution (but with different parameter values and also depending on the likelihood function if the form of the likelihood function is varied). This essentially reduces the difficulty involved in the calculation of the numerical integrations as described above, as the conjugate priors [32]–[34] can provide a closed-form expression for the posterior distribution. Moreover, the use of conjugate priors does not prevent the proposed EM algorithm from choosing flexible forms of the density functions, such as Gaussian, Laplacian, Gamma, or other members in the exponential family, which covers a wide range of distributions for the mixing matrix and the noise covariance. The choice of the hyper-parameters of the priors will be discussed later in Section VI-B. Next, we discuss the choice of the prior densities for the mixing matrix and noise covariance matrix, as in [30], and for the source models, as in [28].

A. Prior Density for Mixing Matrix

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

$$g(a_{ij}|\mu_{ij}, \sigma_{ij}^2) = \mathcal{N}(\mu_{ij}, \sigma_{ij}^2) \quad (5)$$

With (5), some constraints can be imposed on the elements of the mixing matrix, i.e. by assigning some known values to the means μ_{ij} with σ_{ij} chosen for small values to reflect the degree of the uncertainty [30]. Assuming that the elements

¹For a likelihood function, a conjugate prior is defined as the prior for which the posteriori and the priori are of the same type of distributions.

of the mixing matrix are independent from each other, it is straightforward to derive that $g(\text{vec}(\mathbf{A})) = \prod_{i=1}^Q \prod_{j=1}^P g(a_{ij})$, where $\text{vec}(\cdot)$ is an operator for obtaining a vector by stacking the columns of a matrix one beneath the other. It is straightforward to get

$$g(\text{vec}(\mathbf{A})) = \mathcal{N}(\boldsymbol{\mu}_A, \boldsymbol{\Lambda}) \quad (6)$$

where $\boldsymbol{\mu}_A = [\mu_{11}, \dots, \mu_{1P}, \dots, \mu_{Q1}, \dots, \mu_{QP}]^T$ and $\boldsymbol{\Lambda}$ is a diagonal matrix whose elements are σ_{ij}^2 .

B. Prior Density for Noise Covariance Matrix

Covariance matrices are symmetric positive semi-definite matrices. To model the prior knowledge about them, Wishart distribution, which is a generalization of the univariate chi-square distribution, is often used. The Wishart distribution is a conjugate density which therefore has another advantage in simplifying the GMM-EM process as described in Section V. Therefore, as in [30], the Wishart density is assigned as the prior density of the noise covariance matrix \mathbf{R}_w , defined as

$$g(\mathbf{R}_w^{-1} | \Sigma_w^{-1}, v_R) \propto |\mathbf{R}_w^{-1}|^{(v_R - Q - 1)/2} \exp \left[-\frac{1}{2} \text{tr}(\Sigma_w \mathbf{R}_w^{-1}) \right] \quad (7)$$

where Σ_w is a $Q \times Q$ positively defined symmetric matrix, v_R is a scalar greater than $Q - 1$, $\text{tr}(\cdot)$ denotes the trace of a squared matrix, and $|\cdot|$ indicates the determinant of a squared matrix.

C. Prior Density for Parameters of the Source Model

The conjugate prior density assignment for the parameters $\Theta = \{\omega_m, \boldsymbol{\mu}_m, \mathbf{C}_m\}_{m=1}^M$ in (3) is more complicated. According to [35], however, we can interpret a finite mixture density as a density associated with a statistical population, denoted as a mixture of M component populations weighted by the coefficients $(\omega_1, \dots, \omega_M)$. Therefore, we can regard $f(\mathbf{s}(t); \Theta)$ as a marginal PDF of the joint PDF of the parameters Θ . More specifically, it can be computed as the product of a multinomial density and multivariate Gaussian densities, which denote the sizes of the populations and the densities of individual components, respectively [35]. If the joint density of the weighting parameters is a multinomial distribution, then a practical candidate for modeling the prior knowledge of these parameters is a conjugate density such as the Dirichlet density

$$g(\omega_1, \dots, \omega_M | \eta_1, \dots, \eta_M) \propto \prod_{m=1}^M \omega_m^{\eta_m - 1} \quad (8)$$

where $\eta_m > 0$ are the hyper-parameters for the Dirichlet density. As for the parameters $(\boldsymbol{\mu}_m, \mathbf{C}_m)$ of the individual Gaussian mixture component

$$g(\boldsymbol{\mu}_m, \mathbf{C}_m^{-1} | \tau_m, \mathbf{u}_m, v_m, \Sigma_m^{-1}) \propto |\mathbf{C}_m^{-1}|^{(v_m - P)/2} \exp \left[-\frac{\tau_m}{2} (\boldsymbol{\mu}_m - \mathbf{u}_m)^T \mathbf{C}_m^{-1} (\boldsymbol{\mu}_m - \mathbf{u}_m) \right] \exp \left[-\frac{1}{2} \text{tr}(\Sigma_m \mathbf{C}_m^{-1}) \right] \quad (9)$$

where $(\tau_m, \mathbf{u}_m, v_m, \Sigma_m)$ are the prior density hyper-parameters such that $v_m > P - 1$, $\tau_m > 0$, \mathbf{u}_m is a vector of dimension P , and Σ_m is a $P \times P$ positive definite matrix.

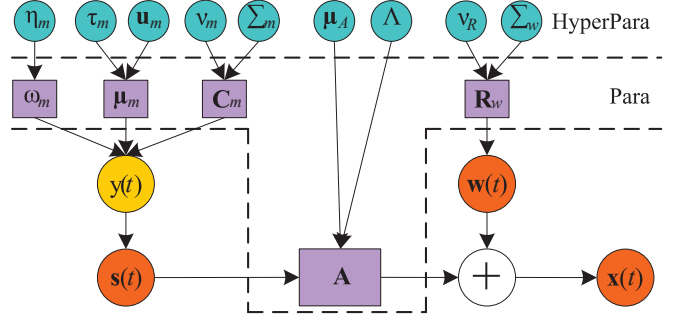


Fig. 1. The probability-generative model of observed signals at the discrete time instance t .

Assuming that the parameters of the individual mixture components and the mixture weights are independent [35], then, the joint prior density $g(\Theta)$ can be computed as the product of the prior densities defined in (8) and (9), respectively, given as follows,

$$g(\Theta) = g(\omega_1, \dots, \omega_M) \prod_{m=1}^M g(\boldsymbol{\mu}_m, \mathbf{C}_m^{-1}) \quad (10)$$

V. BAYESIAN BLIND SEPARATION

In this section, equipped with the source model discussed in Section III and prior densities defined in Section IV, we develop a new EM algorithm under the Bayesian framework for the BSS problem as described in Section II. For the convenience of analysis, we employ a probability-generative model, as depicted in Fig. 1, where a graphical model is used to show the process of generating an observation signal at time instant t based on the mixture model. Apparently, there are two levels of hidden variables in this graphical model, with the first level being represented by the Gaussian component labels $\{y(t)\}_{t=1, \dots, T}$ of the density mixture, and the second level by the source signals $\{\mathbf{s}(t)\}_{t=1, \dots, T}$.

As a result, the BSS problem in essence can be treated as a problem of estimating parameters from incomplete-data. The incomplete-data are the observations $\mathbf{X} = \{\mathbf{x}(t)\}_{t=1, \dots, T}$, while the missing data are the sources $\mathbf{S} = \{\mathbf{s}(t)\}_{t=1, \dots, T}$ and the unobserved Gaussian component labels of the density mixture $Y = \{y(t)\}_{t=1, \dots, T}$. The parameters that need to be estimated are \mathbf{A} , \mathbf{R}_w and Θ . The EM algorithm is a commonly used method for inferring the parameters of an underlying distribution from incomplete data based on the ML/MAP scheme [36]. Therefore, similar to the method we adopted in [30], we derive an GMM-EM algorithm in this work to obtain the MAP estimates of the unknowns including the mixing matrix, noise covariance matrix, and the parameters of the source model, as detailed below.

A. The E-step

Given the observed data \mathbf{X} and the current parameter estimates, the E-step of the GMM-EM algorithm aims to obtain the expected value of the complete-data log-likelihood $\log f(\mathbf{X}, \mathbf{S}, Y | \mathbf{A}, \mathbf{R}_w, \Theta)$ with respect to the unknown data \mathbf{S} and Y . The evaluation of this expectation is called the E-step

of the algorithm. To this end, we define an auxiliary function as

$$J(\mathbf{A}, \mathbf{R}_w, \Theta, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) = \mathbb{E} [\log f(\mathbf{X}, \mathbf{S}, Y | \mathbf{A}, \mathbf{R}_w, \Theta) | \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g] \quad (11)$$

where \mathbf{A}^g , \mathbf{R}_w^g , and Θ^g are the current estimates of the parameters that we use to evaluate the expectation and \mathbf{A} , \mathbf{R}_w , and Θ are the new parameters that we optimize to increase J . $\mathbb{E}[\cdot]$ is an expectation operator.

Since \mathbf{X} and \mathbf{A}^g , \mathbf{R}_w^g , Θ^g are constants, \mathbf{A} , \mathbf{R}_w , Θ are variables that we wish to adjust, and \mathbf{S} , Y are random variables governed by the distribution $f(\mathbf{S}, Y | \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g)$, the RHS of (11) can be rewritten as

$$\mathbb{E} [\log f(\mathbf{X}, \mathbf{S}, Y | \mathbf{A}, \mathbf{R}_w, \Theta) | \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g] = \mathbb{E} [f(\mathbf{S}, Y | \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) \log f(\mathbf{X}, \mathbf{S}, Y | \mathbf{A}, \mathbf{R}_w, \Theta)] \quad (12)$$

After a series of derivations (more details can be found in Appendix A), we can get

$$J = \sum_{t=1}^T \sum_{m=1}^M \int_{\mathbf{s}} f(y(t) = m | \mathbf{s}(t), \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) \log \omega_m f(\mathbf{s}(t) | y(t) = m, \Theta) ds + \sum_{t=1}^T \int_{\mathbf{s}} f(\mathbf{s}(t) | \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) \log f(\mathbf{x}(t) | \mathbf{s}(t), \mathbf{A}, \mathbf{R}_w) ds \quad (13)$$

It is clear that the posterior distribution $f(\mathbf{s}(t) | \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g)$ is indispensable for the evaluation of the expectation in (13). In practice, it can be proved that (more details can be found in Appendix B)

$$f(\mathbf{s}(t) | \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) = \sum_{m=1}^M \tilde{\omega}_{mt}^g \mathcal{N}(\mathbf{s}(t); \tilde{\boldsymbol{\mu}}_{mt}^g, \tilde{\mathbf{C}}_{mt}^g) \quad (14)$$

where

$$\left\{ \begin{array}{l} \tilde{\mathbf{C}}_{mt}^g = ((\mathbf{A}^g)^T (\mathbf{R}_w^g)^{-1} \mathbf{A}^g + (\mathbf{C}_m^g)^{-1})^{-1} \\ \tilde{\boldsymbol{\mu}}_{mt}^g = (\tilde{\mathbf{C}}_{mt}^g) ((\mathbf{A}^g)^T (\mathbf{R}_w^g)^{-1} \mathbf{x}(t) + (\mathbf{C}_m^g)^{-1} \boldsymbol{\mu}_m^g) \\ \tilde{\omega}_{mt}^g = \omega_m^g \left(\frac{|\tilde{\mathbf{C}}_{mt}^g|^{1/2}}{|2\pi \mathbf{R}_w^g|^{1/2} |\mathbf{C}_m^g|^{1/2}} \right) \exp \left\{ -\frac{1}{2} [\mathbf{x}^T(t) (\mathbf{R}_w^g)^{-1} \mathbf{x}(t) + (\boldsymbol{\mu}_m^g)^T (\mathbf{C}_m^g)^{-1} \boldsymbol{\mu}_m^g - (\tilde{\boldsymbol{\mu}}_{mt}^g)^T (\tilde{\mathbf{C}}_{mt}^g)^{-1} \tilde{\boldsymbol{\mu}}_{mt}^g] \right\} \end{array} \right\} \quad (15)$$

B. The M-step

This step maximizes the expectation of the complete-data log-likelihood as shown in (13) with respect to \mathbf{A} , \mathbf{R}_w and Θ , and the maximum point is then taken as the new parameters. It can be observed that the first term in the RHS of (13) is dependent on the parameters Θ of the GMM model, while the second term is determined by the mixing matrix \mathbf{A} and the noise covariance matrix \mathbf{R}_w . For this reason, the auxiliary function in the RHS of (13) can be split into two parts J_1 and J_2 , i.e.

$$J = J_1 + J_2 \quad (16)$$

where

$$\left\{ \begin{array}{l} J_1 = \sum_{t=1}^T \int_{\mathbf{s}} f(\mathbf{s}(t) | \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) \log f(\mathbf{x}(t) | \mathbf{s}(t), \mathbf{A}, \mathbf{R}_w) ds \\ J_2 = \sum_{t=1}^T \sum_{m=1}^M \int_{\mathbf{s}} f(y(t) = m | \mathbf{s}(t), \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) \log \omega_m f(\mathbf{s}(t) | y(t) = m, \Theta) ds \end{array} \right. \quad (17)$$

Accordingly, the parameters \mathbf{A} , \mathbf{R}_w and Θ can be estimated by optimizing the auxiliary functions J_1 and J_2 , respectively, as explained in the next two sub-sections.

1) *Estimation Formula for the Mixing Matrix \mathbf{A} and the Noise Covariance Matrix \mathbf{R}_w* : First of all, the updating rules for the mixing matrix \mathbf{A} and the noise covariance matrix \mathbf{R}_w are discussed based on the auxiliary function J_1 . Note that the auxiliary function J_1 in (17) can be converted into the following form as

$$J_1 = -\frac{T}{2} \log |2\pi \mathbf{R}_w| - \frac{T}{2} \text{tr} [\mathbf{R}_w^{-1} (\mathbf{R}_{xx} - \mathbf{A} \mathbf{R}_{sx} - \mathbf{R}_{sx}^T \mathbf{A}^T + \mathbf{A} \mathbf{R}_{ss} \mathbf{A}^T)] \quad (18)$$

where

$$\left\{ \begin{array}{l} \mathbf{R}_{xx} = \frac{1}{T} \sum_{t=1}^T \mathbf{x}(t) \mathbf{x}^T(t) \\ \mathbf{R}_{sx} = \frac{1}{T} \sum_{t=1}^T \mathbb{E} [\mathbf{s}(t) | \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g] \mathbf{x}^T(t) \\ \mathbf{R}_{ss} = \frac{1}{T} \sum_{t=1}^T \mathbb{E} [\mathbf{s}(t) \mathbf{s}^T(t) | \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g] \end{array} \right. \quad (19)$$

Note that the prior information for the mixing matrix \mathbf{A} denoted by J_A^g is related to its conjugate prior as shown in (6). Therefore, we can obtain the MAP auxiliary function for the mixing matrix \mathbf{A} by incorporating its prior information, which gives

$$\begin{aligned} \hat{J}_A &= J_1 + J_A^g \\ &= -\frac{T}{2} \log |2\pi \mathbf{R}_w| - \frac{1}{2} \log (|\boldsymbol{\Lambda}|) \\ &\quad - \frac{T}{2} \text{tr} [\mathbf{R}_w^{-1} (\mathbf{R}_{xx} - \mathbf{A} \mathbf{R}_{sx} - \mathbf{R}_{sx}^T \mathbf{A}^T + \mathbf{A} \mathbf{R}_{ss} \mathbf{A}^T)] \\ &\quad - \frac{1}{2} \text{tr} [\boldsymbol{\Lambda}^{-1} (\text{vec}(\mathbf{A}) - \boldsymbol{\mu}_A) (\text{vec}(\mathbf{A}) - \boldsymbol{\mu}_A)^T] \end{aligned} \quad (20)$$

The updating rule for the mixing matrix \mathbf{A} can therefore be obtained by taking the derivative of \hat{J}_A with respect to \mathbf{A} , and setting it to zero, which gives

$$\text{vec}(\hat{\mathbf{A}}) = [\mathbf{T} \mathbf{R}_{ss} \otimes \mathbf{R}_w^{-1} + \boldsymbol{\Lambda}^{-1}]^{-1} [\text{vec}(\mathbf{T} \mathbf{R}_w^{-1} \mathbf{R}_{xs}) + \boldsymbol{\Lambda}^{-1} \boldsymbol{\mu}_A] \quad (21)$$

where \otimes denotes the Kronecker product.

Similarly, by incorporating the prior information for the noise covariance matrix \mathbf{R}_w denoted by $J_{R_w}^g$, which is related to its prior distribution as defined in (7), the MAP auxiliary

function for the noise covariance \mathbf{R}_w can be written as

$$\begin{aligned}\hat{J}_{R_w} &= J_1 + J_{R_w}^g \\ &= -\frac{T}{2} \log |\mathbf{R}_w| \\ &\quad - \frac{T}{2} \text{tr} [\mathbf{R}_w^{-1} (\mathbf{R}_{xx} - \mathbf{A} \mathbf{R}_{sx} - \mathbf{R}_{sx}^T \mathbf{A}^T + \mathbf{A} \mathbf{R}_{ss} \mathbf{A}^T)] \\ &\quad - \frac{v_R - Q - 1}{2} \log |\mathbf{R}_w| - \frac{1}{2} \text{tr} (\Sigma_w \mathbf{R}_w^{-1})\end{aligned}\quad (22)$$

The updating rule for the noise covariance matrix \mathbf{R}_w can be similarly obtained by taking the derivative of \hat{J}_{R_w} with respect to \mathbf{R}_w , and setting it to zero, which gives

$$\hat{\mathbf{R}}_w = \frac{1}{T + v_R - Q - 1} \left[T (\mathbf{R}_{xx} - \hat{\mathbf{A}} \mathbf{R}_{sx} - \mathbf{R}_{sx}^T (\hat{\mathbf{A}})^T + \hat{\mathbf{A}} \mathbf{R}_{ss} (\hat{\mathbf{A}})^T) + \Sigma_w \right] \quad (23)$$

The updating rules for \mathbf{A} and \mathbf{R}_w involve the calculation of \mathbf{R}_{ss} and \mathbf{R}_{sx} . Using the posterior distribution $f(\mathbf{s}(t) | \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g)$ as shown in (14) and (15), it is easy to obtain the conditional expectations $\mathbb{E}[\mathbf{s}(t) | \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g]$ and $\mathbb{E}[\mathbf{s}^T(t) \mathbf{s}(t) | \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g]$.

2) *Estimation Formula for the GMM Parameters:* The parameters Θ of the GMM can be updated in a similar way to that for the mixing matrix and noise covariance matrix. Using (14) and (15), the auxiliary function J_2 in (17) can be re-written as

$$J_2 = \sum_{t=1}^T \sum_{m=1}^M \int_{\mathbf{s}} \tilde{\omega}_{mt}^g \mathcal{N}(\mathbf{s}(t); \tilde{\boldsymbol{\mu}}_{mt}^g, \tilde{\mathbf{C}}_{mt}^g) \log \omega_m f(\mathbf{s}(t) | y(t) = m, \Theta) ds \quad (24)$$

with the prior density as depicted in (8), (9) and (10), then the prior information for the source signals can be denoted as

$$\begin{aligned}J_S^g &= \sum_{m=1}^M (\eta_m - 1) \log \omega_m - ((v_m - P)/2) \log |\mathbf{C}_m| \\ &\quad - \frac{\tau_m}{2} \text{tr} [\mathbf{C}_m^{-1} (\boldsymbol{\mu}_m - \mathbf{u}_m) (\boldsymbol{\mu}_m - \mathbf{u}_m)^T] - \frac{1}{2} \text{tr} (\Sigma_m \mathbf{C}_m^{-1})\end{aligned}\quad (25)$$

Hence, the MAP auxiliary function \hat{J}_S for the GMM parameters can be written as $\hat{J}_S = J_2 + J_S^g$. To maximize this expression, we can maximize the terms containing the weighting coefficient parameter ω_m , and the term containing the mean vector $\boldsymbol{\mu}_m$, and the covariance matrix \mathbf{C}_m , $m = 1, \dots, M$ separately since they are independent from each other.

Note that $\int_{\mathbf{s}} \mathcal{N}(\mathbf{s}(t); \tilde{\boldsymbol{\mu}}_{mt}^g, \tilde{\mathbf{C}}_{mt}^g) ds = 1$, hence, the part of the auxiliary function \hat{J}_S related to parameter ω_m can be simplified as

$$\hat{J}_S^{(1)} = \sum_{t=1}^T \sum_{m=1}^M \tilde{\omega}_{mt}^g \log \omega_m + \sum_{m=1}^M (\eta_m - 1) \log \omega_m \quad (26)$$

Adding the Lagrange multiplier λ , using the constraints that $\sum_{m=1}^M \omega_m = 1$, and setting the derivative of $\hat{J}_S^{(1)}$ with respect

to ω_m equal to zero, one obtains

$$\begin{aligned}\frac{\partial}{\partial \omega_m} \left[\left(\sum_{m=1}^M \sum_{t=1}^T \tilde{\omega}_{mt}^g \log \omega_m \right) \right. \\ \left. + \sum_{m=1}^M (\eta_m - 1) \log \omega_m + \lambda \left(\sum_{m=1}^M \omega_m - 1 \right) \right] = 0\end{aligned}\quad (27)$$

Summing both sides over m , we can get $\lambda = -(\sum_{m=1}^M \eta_m - M + T)$ resulting in

$$\hat{\omega}_m = \frac{\eta_m - 1 + \sum_{t=1}^T \tilde{\omega}_{mt}^g}{\sum_{m=1}^M \eta_m - M + T} \quad (28)$$

On the other hand, the part of the auxiliary function \hat{J}_S related to parameters $\boldsymbol{\mu}_m$ and \mathbf{C}_m can be written as

$$\begin{aligned}\hat{J}_S^{(2)} &= \sum_{t=1}^T \int_{\mathbf{s}} \tilde{\omega}_{mt}^g \mathcal{N}(\mathbf{s}(t); \tilde{\boldsymbol{\mu}}_{mt}^g, \tilde{\mathbf{C}}_{mt}^g) \log f(\mathbf{s}(t) | y(t) = m, \Theta) ds \\ &\quad + ((v_m - P)/2) \log |\mathbf{C}_m^{-1}| \\ &\quad - \frac{\tau_m}{2} \text{tr} [\mathbf{C}_m^{-1} (\boldsymbol{\mu}_m - \mathbf{u}_m) (\boldsymbol{\mu}_m - \mathbf{u}_m)^T] - \frac{1}{2} \text{tr} [\Sigma_m \mathbf{C}_m^{-1}]\end{aligned}\quad (29)$$

The updating rule for $\boldsymbol{\mu}_m$ and \mathbf{C}_m can therefore be obtained by taking the derivative of $\hat{J}_S^{(2)}$ with respect to $\boldsymbol{\mu}_m$ and \mathbf{C}_m respectively, and setting them to zero. That is $\partial \hat{J}_S^{(2)} / \partial \boldsymbol{\mu}_m = 0$ and $\partial \hat{J}_S^{(2)} / \partial \mathbf{C}_m = 0$.

For notational simplicity, define

$$\begin{aligned}\boldsymbol{\Gamma}_m &= \mathbb{E}[\mathbf{s}(t) \mathbf{s}^T(t) | y(t) = m, \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g] \\ &= \tilde{\mathbf{C}}_{mt}^g + \tilde{\boldsymbol{\mu}}_{mt}^g (\tilde{\boldsymbol{\mu}}_{mt}^g)^T\end{aligned}\quad (30)$$

Then, one obtains

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\mu}_m} \left(\sum_{t=1}^T \int_{\mathbf{s}} \tilde{\omega}_{mt}^g \mathcal{N}(\mathbf{s}(t); \tilde{\boldsymbol{\mu}}_{mt}^g, \tilde{\mathbf{C}}_{mt}^g) \right. \\ \left. \log f(\mathbf{s}(t) | y(t) = m, \Theta) ds \right) \\ = \sum_{t=1}^T \tilde{\omega}_{mt}^g \mathbf{C}_m^{-1} (\tilde{\boldsymbol{\mu}}_{mt}^g - \boldsymbol{\mu}_m)\end{aligned}\quad (31)$$

and

$$\begin{aligned}\frac{\partial}{\partial \mathbf{C}_m} \left(\sum_{t=1}^T \int_{\mathbf{s}} \tilde{\omega}_{mt}^g \mathcal{N}(\mathbf{s}(t); \tilde{\boldsymbol{\mu}}_{mt}^g, \tilde{\mathbf{C}}_{mt}^g) \right. \\ \left. \log f(\mathbf{s}(t) | y(t) = m, \Theta) ds \right) \\ = \frac{1}{2} \sum_{t=1}^T \tilde{\omega}_{mt}^g [-\mathbf{C}_m + (\boldsymbol{\Gamma}_m - \tilde{\boldsymbol{\mu}}_{mt}^g \boldsymbol{\mu}_m^T - \boldsymbol{\mu}_m (\tilde{\boldsymbol{\mu}}_{mt}^g)^T + \boldsymbol{\mu}_m \boldsymbol{\mu}_m^T)]\end{aligned}\quad (32)$$

On the other hand, notice that

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\mu}_m} \left(((v_m - P)/2) \log |\mathbf{C}_m^{-1}| - \right. \\ \left. \frac{\tau_m}{2} \text{tr} [\mathbf{C}_m^{-1} (\boldsymbol{\mu}_m - \mathbf{u}_m) (\boldsymbol{\mu}_m - \mathbf{u}_m)^T] - \frac{1}{2} \text{tr} (\Sigma_m \mathbf{C}_m^{-1}) \right) \\ = -\tau_m \mathbf{C}_m^{-1} (\boldsymbol{\mu}_m - \mathbf{u}_m)\end{aligned}\quad (33)$$

$$\begin{aligned} & \frac{\partial}{\partial \mathbf{C}_m} \left(((v_m - P)/2) \log |\mathbf{C}_m^{-1}| - \right. \\ & \quad \left. \frac{\tau_m}{2} \text{tr}[\mathbf{C}_m^{-1}(\boldsymbol{\mu}_m - \mathbf{u}_m)(\boldsymbol{\mu}_m - \mathbf{u}_m)^T] - \frac{1}{2} \text{tr}(\Sigma_m \mathbf{C}_m^{-1}) \right) \\ & = -\frac{(v_m - P)}{2} \mathbf{C}_m + \frac{\tau_m}{2} (\boldsymbol{\mu}_m - \mathbf{u}_m)(\boldsymbol{\mu}_m - \mathbf{u}_m)^T + \frac{1}{2} \Sigma_m \end{aligned} \quad (34)$$

Combining the terms (31) and (33), the updating formula for $\boldsymbol{\mu}_m$ can be easily obtained as

$$\hat{\boldsymbol{\mu}}_m = \frac{\tau_m \mathbf{u}_m + \sum_{t=1}^T \tilde{\omega}_{mt}^g \tilde{\boldsymbol{\mu}}_{mt}^g}{\tau_m + \sum_{t=1}^T \tilde{\omega}_{mt}^g} \quad (35)$$

and the updating formula for \mathbf{C}_m can be similarly obtained by combining the terms (32) and (34)

$$\begin{aligned} \hat{\mathbf{C}}_m = & \frac{\Sigma_m + \tau_m (\hat{\boldsymbol{\mu}}_m - \mathbf{u}_m)(\hat{\boldsymbol{\mu}}_m - \mathbf{u}_m)^T +}{v_m - P + \sum_{t=1}^T \tilde{\omega}_{mt}^g} + \\ & \frac{\sum_{t=1}^T \tilde{\omega}_{mt}^g (\mathbf{\Gamma}_m - \tilde{\boldsymbol{\mu}}_{mt}^g (\hat{\boldsymbol{\mu}}_m)^T - \hat{\boldsymbol{\mu}}_m (\tilde{\boldsymbol{\mu}}_{mt}^g)^T + \hat{\boldsymbol{\mu}}_m (\hat{\boldsymbol{\mu}}_m)^T)}{v_m - P + \sum_{t=1}^T \tilde{\omega}_{mt}^g} \end{aligned} \quad (36)$$

VI. PRACTICAL IMPLEMENTATION AND ALGORITHM ANALYSIS

In this section, we discuss some practical implementation issues of the proposed algorithm, and also offer an empirical analysis of its convergence and computational complexity.

A. Summary of the Algorithm

The proposed GMM-EM algorithm can be implemented as follows:

- 1) **Initialization** Initialize the mixing matrix $\mathbf{A}^{(0)}$, noise variance matrix $\mathbf{R}_w^{(0)}$ and model parameters $\Theta^{(0)}$ according to the initialization scheme described in Section VI-B, and set the EM iteration index $i = 0$.
- 2) **EM iterations** Repeat the E-step and M-step until convergence.
 - a) **E-step** Calculate $f(\mathbf{s}(t)|\mathbf{x}(t), \mathbf{A}^{(i)}, \mathbf{R}_w^{(i)}, \Theta^{(i)})$ according to (14)-(15), and calculate \mathbf{R}_{ss} and \mathbf{R}_{sx} according to (19).
 - b) **M-step** Calculate the mixing matrix $\mathbf{A}^{(i+1)}$ and noise co-variance matrix $\mathbf{R}_w^{(i+1)}$ according to (21) and (23), respectively. Calculate the weight $\omega_m^{(i+1)}$, mean vector $\boldsymbol{\mu}_m^{(i+1)}$ and co-variance matrix $\mathbf{C}_m^{(i+1)}$ according to (28), (35) and (36), respectively.
- 3) **MAP source estimation** Let i_0 be the number of iterations required before the convergence of the algorithm, the posterior mean estimate $\hat{\mathbf{s}}_{MAP}$ is approximated by the empirical mean of the sequence $\mathbf{s}_{i>i_0}$.

B. Parameter Initialization

It is well-known that the EM optimization strategy is sensitive to the initial set-ups of the parameters. The likelihood function may converge to a local maximum instead of the global maximum due to the use of the bootstrap process in

the iterations. The initial values for the parameters therefore become important for the convergence of the EM algorithm, not only in terms of minimizing the number of iterations required for the algorithm to converge to a local maximum, but also for it to find a ‘‘good’’ solution [33]. Therefore, the parameters are given proper prior densities, i.e., conjugate prior densities, in order to incorporate the prior information as discussed in Section IV. A reasonable choice for the initial estimates is the mode of the prior density. For the mixing matrix and the noise covariance matrix, the initial values can be set as

$$\begin{cases} \text{vec}(\mathbf{A})^{(0)} = \boldsymbol{\mu}_A \\ \mathbf{R}_w^{(0)} = (v_R - Q - 1) \Sigma_w^{-1} \end{cases} \quad (37)$$

The mean value of the mixing matrix $\boldsymbol{\mu}_A$ can be estimated by the blind identification methods such as those presented in [6]–[12]. Although, in such cases, the source signals can not be recovered by multiplying the observed signals with the inverse/pseudo inverse of the mixing matrix, the hyperparameter $\boldsymbol{\mu}_A$ can be determined by the estimate of the mixing matrix. In our implementation in Section VII, the LEMACAF-4 method² in [10] has been used for estimating the initial value of the mixing matrix.

Similarly, the initial estimates for the GMM model parameters of the sources are taken as

$$\begin{cases} \omega_m^{(0)} = (\eta_m - 1) / \left(\sum_{m=1}^M \eta_m - M \right) \\ \boldsymbol{\mu}_m^{(0)} = \mathbf{u}_m \\ \mathbf{C}_m^{(0)} = (v_m - P) \Sigma_m^{-1} \end{cases} \quad (38)$$

It has been shown that the joint PDF of the observed signals can also be modeled by GMM when the joint PDF of the source signals is modeled by GMM [28] (due to space limitation, the detail is omitted here). As a result, the weighting coefficients, the mean vectors and covariance matrices of the observation-based GMM model can be learned from the observed signals. Therefore, according to the relationship between the weighting coefficients, the mean vectors and covariance matrices of the observation-based GMM model and their counterparts of the source-based GMM model, the hyperparameters η_m , \mathbf{u}_m and Σ_m can be determined by the estimate of the mixing matrix and the estimates of the observation-based GMM model parameters jointly.

C. Convergence Analysis

In essence, the proposed GMM-EM separating algorithm is a gradient-based bootstrap method for optimizing the log-likelihood function. There are already some works that have investigated this issue [37], [38]. For example, Xu *et al.* [37] have established the linkage of the EM algorithm with the gradient-based approaches for the ML learning based on GMM, and shown that the EM parameters are iterated in terms of the gradient obtained by a positive definite projection matrix. This result was extended to the more general block

²Matlab codes can be found at: <http://www.i3s.unice.fr/pcomon/TensorPac/kage.html>

coordinate descent (BCD) method [39], [40], in which a single block of variables is optimized at each iteration.

In Section V, to update each variable of \mathbf{A} , \mathbf{R}_w , ω_m , $\boldsymbol{\mu}_m$ and \mathbf{C}_m , the employed method is to simply set the first derivative of the complete-data log-likelihood with respect to each variable to be zero and solve the corresponding equation sequentially. Hence, if the Hessian matrices of complete-data log-likelihood with respect to these variables are non-positive, then it is safe to state that the subproblem with respect to each variable is convex. Take the mixing matrix as an example. By converting the matrix \mathbf{A} into vector form $\text{vec}(\mathbf{A})$, the Hessian matrix of \hat{J}_A with respect to $\text{vec}(\mathbf{A})$ can be easily obtained, and written as $\mathbf{H}_A = -(T\mathbf{R}_{s_s} \otimes \mathbf{R}_w^{-1} + \boldsymbol{\Lambda}^{-1})$. Since \mathbf{R}_w and $\boldsymbol{\Lambda}$ are positively defined, it is obvious that \mathbf{H}_A is non-positive. The Hessian matrix of \hat{J}_{R_w} with respect to \mathbf{R}_w can be obtained in a similar way, and it is also non-positive. That is, for each iteration of the GMM-EM algorithm given in (21, 23), the search direction of the parameters has a positive projection on the gradient of its corresponding MAP auxiliary function.

We further discuss the GMM parameters ω_m , $\boldsymbol{\mu}_m$ and \mathbf{C}_m . If each mixture component is assumed to be non-degenerate [37], i.e., $\hat{\omega}_m > 0$, then $\tilde{\omega}_{m1}^g, \dots, \tilde{\omega}_{mT}^g$ is a sequence of T i.i.d. random variables with a non-degenerate distribution and $\lim_{T \rightarrow \infty} \sum_{t=1}^T \tilde{\omega}_{mt} = \infty$ with probability one. It follows that the Hessian matrix $\mathbf{H}_{\omega_m} = -(\eta_m - 1 + \sum_{t=1}^T \tilde{\omega}_{mt}^g)$ is non-positive with probability one when $T \rightarrow \infty$. Applying the same reasoning, we can see that the GMM-EM estimation formulas for $\hat{\boldsymbol{\mu}}_m$ and $\hat{\mathbf{C}}_m$ are asymptotically similar in terms of the MAP approach [35], [41]. Therefore, as long as the initial estimates of $\mathbf{A}^{(0)}$, $\mathbf{R}_w^{(0)}$ and $\Theta^{(0)}$ remain unchanged, the EM algorithm will converge to the same estimates with probability one when $T \rightarrow \infty$.

Finally, it should be pointed out that the EM algorithm may converge to a local maximum instead of the global maximum when the number of parameters is large and/or the parameters of the algorithm are inappropriately initialized. This is a general limitation associated with the gradient-based bootstrap-like optimization algorithms. The reason is that it is often trapped to the neighborhood of a local optimizer if the number of parameters is large and/or the parameters of the EM algorithm are initialized such that the solution is far from the global optimizer. Hence, the initialization scheme discussed in Section VI-B is vital to ensure the convergence of the proposed GMM-EM algorithm.

D. Computational Complexity Analysis

The computational load of the proposed GMM-EM algorithm is dominated by the E-step and M-step. In each iteration of the E-step, it is required to:

- calculate the posterior probability $f(s(t)|\mathbf{X}, \mathbf{A}^{(i)}, \mathbf{R}_w^{(i)}, \Theta^{(i)})$ with (14) and (15) which requires $\mathcal{O}(Q(P+Q)(P^2+Q^2))$ multiplications per observation vector.
- calculate the statistics \mathbf{R}_{xx} , \mathbf{R}_{sx} and \mathbf{R}_{ss} with (19) which requires $\mathcal{O}((P+Q)^2)$ multiplications per observation vector.

In each iteration of the M-step, it is required to:

- update the mixing matrix \mathbf{A} using (21) and the noise covariance matrix \mathbf{R}_w using (23) which require $\mathcal{O}((PQ)^3 + (PQ)^2)$ and $\mathcal{O}(PQ(P+Q))$ multiplications, respectively.
- update the weighting coefficients $\omega_m^{(i+1)}$, mean vector $\boldsymbol{\mu}_m^{(i+1)}$ and covariance matrix $\mathbf{C}_m^{(i+1)}$ according to (28), (35) and (36) which amount to $\mathcal{O}(MP^2T)$ multiplications.

From the above analysis, we can see that the computational complexity of the proposed GMM-EM algorithm depends closely on the number of sources and sample size considered in the model. Theoretically, the proposed GMM-EM algorithm for GMM parameter estimation of source signals would become increasingly intractable and computationally unaffordable as the number of sources increases. This is because the number of Gaussians for modeling the source vector grows exponentially with the number of sources. For example, assuming that the number of sources is $P = 10$, and the PDF of each source is modeled by the GMM with $l_p = 3$ Gaussians, then it is straightforward to derive that $M = 3^{10}$ Gaussians are required to model the source vector. However, it has been shown that the determined GMM order in high dimensions is always much smaller than the theoretical number of Gaussians [28]. The main reason is that the distribution of the sensors becomes more Gaussian while the number of sources increases. Hence, it enables the applicability of the proposed GMM-EM method also for a large number of sources.

Note that the computational complexity of an iteration of the NCA algorithm [14] is $\mathcal{O}(T^3)$ when the size of the signal is assumed to be much larger than other parameters. In contrast to the NCA algorithm, the proposed GMM-EM algorithm apparently has advantage in terms of computational complexity. This is because the computational complexity of an iteration of the proposed GMM-EM algorithm is $\mathcal{O}(T)$ under the same situation.

VII. SIMULATIONS AND ANALYSIS

In this section, the separation performance of the proposed GMM-EM algorithm is evaluated in terms of similarity score, and compared with that of the NCA [14] and FastICA [4] algorithms. Calculation of the similarity score is detailed in Appendix C. Note that the FastICA algorithm cannot be implemented in the underdetermined case, hence, we only compare the proposed GMM-EM with the NCA in such case.

The section is organized as follows. Firstly, the separation performances of the compared algorithms are evaluated based on synthetic data in terms of similarity score versus the signal-to-noise (SNR) level within the mixtures, sample size, and the number of sources in determined mixtures. Secondly, these performance aspects are also investigated for underdetermined mixtures. Finally, the performances of the compared algorithms are evaluated for separating mixtures of real speech signals.

The compared algorithms were operated under the following overall settings: 1) The number of EM iterations used in the proposed GMM-EM algorithm was set to 100; 2) The separation performance of the NCA algorithm was evaluated with 100 iterations.

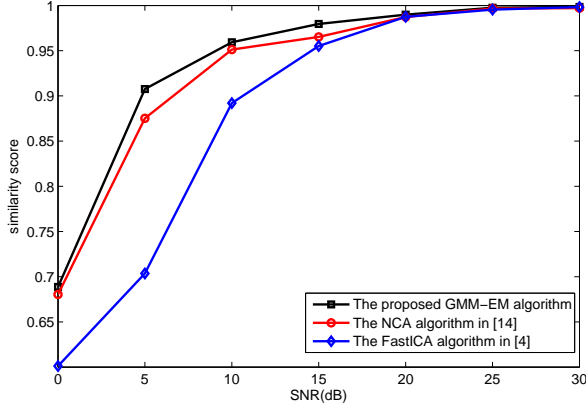


Fig. 2. The average similarity score of the tested algorithms versus the SNR in the determined case.

A. Synthetic Data

1) *Separation Performance as a Function of SNR in Determined Mixtures:* The following experiment compares the separation performances of the tested algorithms for determined mixtures in the presence of additive Gaussian noise. Each source signal is synthesized by the following GMM PDF $f_s = 0.5\mathcal{N}(s; 1, 0.4) + 0.5\mathcal{N}(s; -1, 0.4)$. For each SNR, 100 sets of two-dimensional independent source signals, containing $T = 1000$ samples, are synthesized and mixed by a random 2×2 mixing matrix, with its elements randomly drawn in the range $[-1, 1]$. Additive Gaussian noise is added in the mixing process, and the SNR of the observations ranges from 0 dB to 30 dB.

Fig. 2 depicts the average similarity score between the original sources and recovered sources of the tested algorithms versus the SNR. For the low SNRs (e.g., lower than 20 dB), one can observe that the proposed GMM-EM algorithm offers the best performance, followed by the NCA and FastICA algorithms respectively. The advantage of the proposed algorithm tends to disappear when the SNR is greater than 20 dB, and in this case, the level of noise is pretty low and hence could be ignored in practice. Furthermore, the performance of the NCA algorithm is close to that of the proposed GMM-EM algorithm. This is because the noise component is also considered in the NCA algorithm. The main reason for the proposed GMM-EM algorithm being robust to noise is that the noise component has been taken into account in the model with its covariance jointly estimated in the EM process.

2) *Separation Performance as a Function of Sample Size in Determined Mixtures:* The following experiment compares the separation performances of the tested algorithms as a function of the sample size T for determined mixtures. Each signal is synthesized by the same GMM used in the first experiment shown above. For each $T \in \{100, 200, 400, 600, 1000, 2000, 4000\}$, 100 sets of two-dimensional independent sources are synthesized and mixed by a random 2×2 mixing matrix, with its elements randomly drawn in the range $[-1, 1]$. Additive Gaussian noise is added in the mixing process, and the SNR of the observations is 10

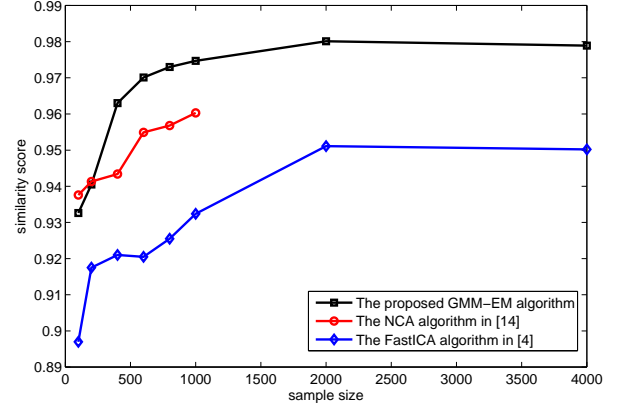


Fig. 3. The average similarity score of the tested algorithms versus the sample size in the determined case.

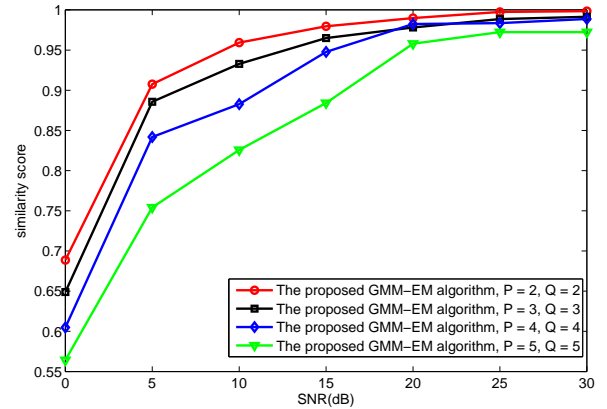


Fig. 4. The average similarity score of the proposed algorithm versus SNR for a varying number of sources in the determined case.

dB.

Fig. 3 depicts the average similarity score of the tested algorithm versus the sample size T . One can observe that as T increases from 100 to 4000, the separation performances of the tested algorithms improve, and the proposed GMM-EM and NCA algorithms outperform the FastICA algorithm. However, as pointed out in [14], the computational complexity of the NCA algorithm is proportional to $\mathcal{O}(T^3)$, hence, it becomes computationally prohibitive when the sample size is large than 1000, and no results are given beyond this point.

3) *Separation Performance as a Function of Dimension in Determined Mixtures:* The following experiment compares the separation performances of the proposed algorithm as a function of SNR for different number of sources in determined mixtures. Each signal is synthesized by the same GMM as used in the first experiment. For each $P \in \{2, 3, 4, 5\}$, 100 sets of P independent source signals, containing $T = 1000$ samples, are synthesized and mixed by a random $P \times P$ mixing matrix, with its elements randomly drawn in the range $[-1, 1]$. Additive Gaussian noise is added in the mixing process, and the SNR of the observations ranges from 0 dB to 30 dB.

Fig. 4 depicts the average similarity score of the proposed

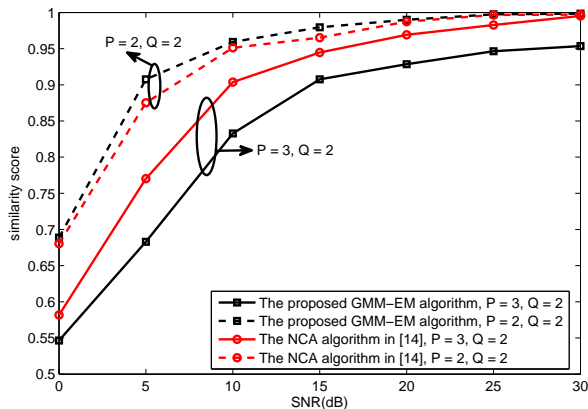


Fig. 5. The average similarity scores of the tested algorithms versus the SNR in the underdetermined case.

GMM-EM algorithm versus the SNR when the number of sources varied from 2 to 5. One can observe that the performance of the proposed GMM-EM algorithm deteriorates as the number of sources increases.

4) *Separation Performance as a Function of SNR in Underdetermined Mixture:* The following experiment compares the separation performances of the tested algorithms for underdetermined mixtures in the presence of additive Gaussian noise. Each signal is synthesized by the same GMM used in the first experiment. For each SNR, 100 sets of three-dimensional independent source signals, containing $T = 1000$ samples, are synthesized and mixed by a random 2×3 mixing matrix, with its elements randomly drawn from the range $[-1, 1]$. Additive Gaussian noise is added in the mixing process, and the SNR of the observations ranges from 0 dB to 30 dB.

The average similarity scores of the tested algorithms versus the SNR in the underdetermined case are shown in Fig. 5. Due to the mixing matrix and sources are estimated jointly, rather than separately (e.g., based on the inverse of the mixing matrix), the proposed GMM-EM algorithm can also work in the underdetermined case. From Fig. 5, one can also observe that the separation performances of the tested algorithms improve with the increase of SNR. However, it can also be observed that the separation performance in the underdetermined case deteriorates as compared with the performance in determined case. This is because the information loss caused by the lack of sensors in the underdetermined case. Furthermore, it should be pointed out that the NCA algorithm outperforms the proposed GMM-EM algorithm in such an underdetermined case. The main reason is that the NCA algorithm mainly depends on whether the null spaces of different sources are orthogonal, regardless whether the mixture is determined or underdetermined.

5) *Separation Performance as a Function of Sample Size in Underdetermined Mixtures:* The following experiment compares the separation performances of the tested algorithms as a function of sample size T for underdetermined mixtures. Each signal is synthesized by the same GMM as used in the first experiment. For each $T \in \{100, 200, 400, 600, 1000, 2000, 4000\}$, 100 sets of three-

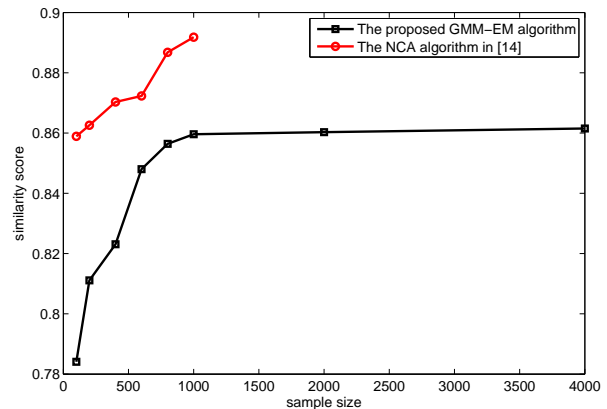


Fig. 6. The average similarity score of the tested algorithms versus the sample size in the underdetermined case.

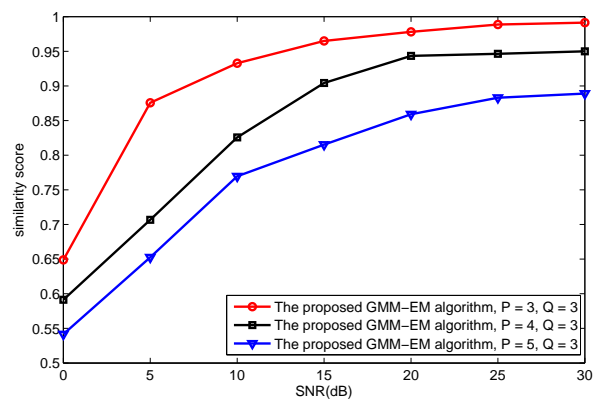


Fig. 7. The average similarity scores of the proposed algorithm versus the SNR for a varying number of sources in the underdetermined case.

dimensional independent source signals are synthesized and mixed by a 2×3 mixing matrix, with its elements randomly drawn in the range $[-1, 1]$. Additive Gaussian noise is added in the mixing process, and the SNR of the observations is 10 dB.

Fig. 6 depicts the average similarity score of the tested algorithm versus the sample size T , where the performance of the NCA algorithm is again shown for up to 1000 samples. One can observe that as T increases from 100 to 4000, the separation performances of the tested algorithms improve.

6) *Separation Performance as a Function of Dimension in Underdetermined Mixtures:* The following experiment compares the separation performances of the tested algorithms as a function of the number of sources in underdetermined mixtures. Each signal is synthesized by the same GMM used in the first experiment. For each $P \in \{3, 4, 5\}$, 100 sets of P independent source signals, containing $T = 1000$ samples, are synthesized and mixed by a random $3 \times P$ mixing matrix, with its elements randomly drawn in the range $[-1, 1]$. Additive Gaussian noise is added in the mixing process, and the SNR of the observations ranges from 0 dB to 30 dB.

Fig. 7 depicts the average similarity score of the proposed GMM-EM algorithm versus the SNR when the number of

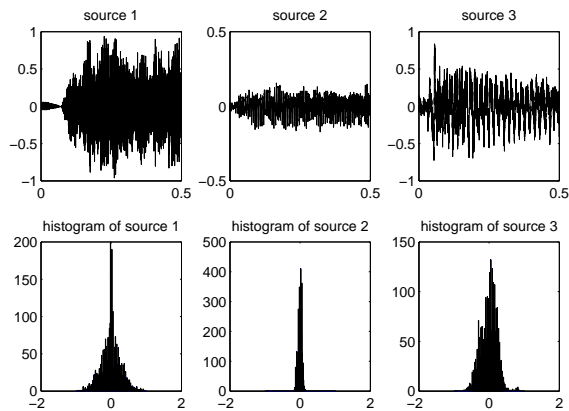


Fig. 8. Speech signals and their respective histograms of amplitude distributions.

observations is fixed and the number of sources is varied. One can observe that the performance of the proposed GMM-EM algorithm deteriorates as the number of sources increases.

B. Real Data

The following experiments compare the performances of the tested algorithms, in terms of similarity score, in separating different mixture combinations of three 0.5-s-long speech signals³, sampled at 8000 Hz and recorded with 8 bits per sample. Fig. 8 shows the waveform of the sources and their histograms of the amplitude distributions. On the other hand, the PDF of each source signal is modeled by GMM of order 3, where the order of the GMM is determined according to the Bayesian information criterion (BIC) [42]. The BIC, based on the likelihood function and a penalty term introduced for the number of parameters in the model, is a well known criterion for model selection among a finite set of models. By calculating the BIC values for all possible models, the candidate model is chosen as the one corresponding to the minimum value of the BIC. The density estimations for the sources are shown in Fig. 9. We can observe that each distribution can be well approximated with 3 Gaussian components and similar to its counterpart in Fig. 8.

Two experiments are used to investigate the separation performance of the proposed GMM-EM algorithm when the sources are real speech signals.

- In the first experiment, two speech signals are artificially mixed by a 2×2 mixing matrix, whose elements are randomly generated from a uniform distribution over the interval $[-1, 1]$. Additive Gaussian noise is added in the mixing process, and the SNR of the observations ranges from 0 dB to 30 dB. In [28], it is shown that the joint PDF of the observed signals can also be modeled by GMM when the joint PDF of the source signals is modeled by GMM. Hence, the order of the GMM of the sources equals to that of the observed signals. As a result, the order of the GMM can be determined

³available at: <http://www.kecl.ntt.co.jp/icl/signal/sawada/webdemo/bssdemo.html>

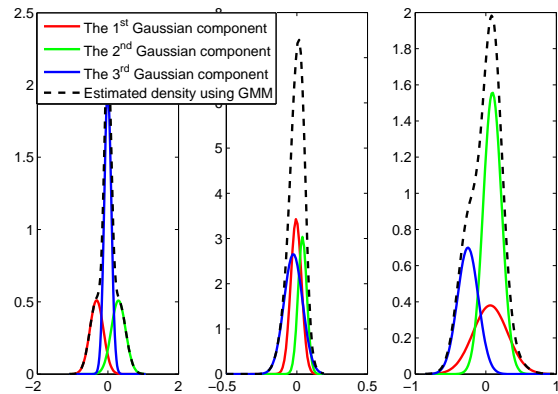


Fig. 9. GMM for probability density estimations of speech signals.

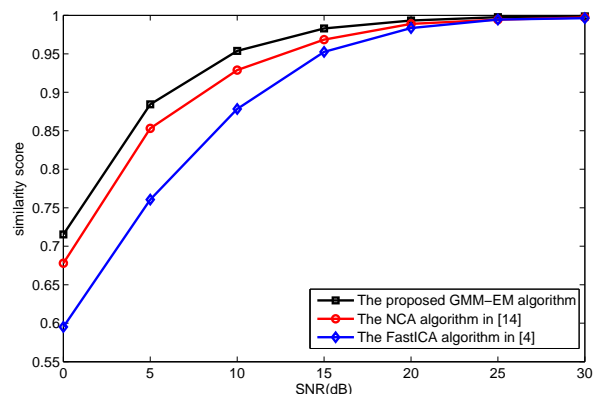


Fig. 10. The average similarity score of the tested algorithms versus the SNR in the determined mixtures with real speech signals.

according to the BIC based on the observed signals. Here, the optimal GMM order determined by the BIC criterion for the GMM-EM separating algorithm is 9. The NCA and FastICA algorithms are implemented as baseline algorithms. 100 Monte Carlo experiments are run.

- In the second experiment, the underdetermined case of $P = 3$ sources and $Q = 2$ observations is considered. The source signals are the same speech signals used in the noisy determined case. The sources are artificially mixed by a 2×3 mixing matrix, whose elements are randomly generated from a uniform distribution over the interval $[-1, 1]$. The SNR of the observations is ranged from 0 dB to 30 dB. The optimal GMM order determined by the BIC criterion for the GMM-EM algorithm is 27. The NCA algorithm is implemented as the baseline algorithm. 100 Monte Carlo experiments are run.

The average separation performance of the tested algorithms in terms of the similarity score versus the SNR in determined mixtures is shown in Fig. 10. The separation performances of the tested algorithms for underdetermined mixtures are shown in Fig. 11. One can observe similar performance patterns to those in Fig. 2 and Fig. 5. More specifically, the proposed

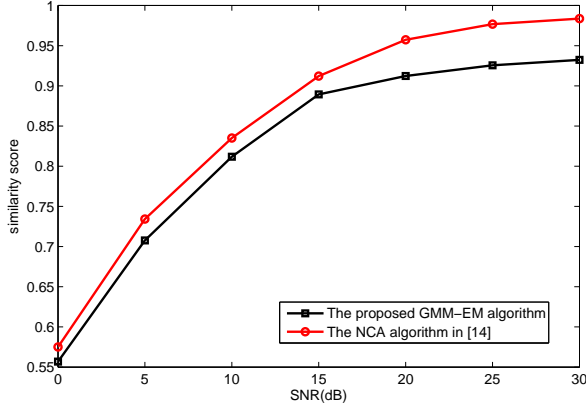


Fig. 11. The average similarity score of the tested algorithms versus the SNR in the underdetermined mixtures when the sources are speech signals.

GMM-EM algorithm seems to give almost identical similarity scores as for synthetic data.

VIII. CONCLUSION

In this paper, the challenging noisy and/or underdetermined BSS problem is considered. To address this issue, we have proposed a GMM-EM approach in which the non-Gaussianity of the sources is exploited by modeling their distributions using GMM. Then, the mixing coefficients, the GMM parameters and the noise covariance matrix are estimated by maximizing their posterior probabilities using an EM algorithm. Finally, issues regarding the practical implementation and performance of the proposed GMM-EM algorithm, such as the initialization scheme for the parameters, the convergence performance and computational complexity, are also discussed. Simulation results have shown that the proposed GMM-EM algorithm gives promising results in two difficult cases: low SNR and underdetermined mixtures. Taking into account the noise in the model and jointly estimating its covariance are the main reasons for the robust performance achieved by the proposed GMM-EM method in noisy environments. The competitive separation performance achieved by the proposed algorithm in underdetermined cases is mainly due to the incorporation of prior information by conjugate priors to assist recovering the sources and no computation for the inverse of the mixing matrix.

APPENDIX A PROOF OF EQUATION (13)

Since

$$\begin{aligned} f(\mathbf{S}, Y | \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) \\ = \prod_{t=1}^T f(\mathbf{s}(t) | \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) f(y(t) | \mathbf{s}(t), \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) \end{aligned} \quad (39)$$

On the other hand,

$$\begin{aligned} f(\mathbf{X}, \mathbf{S}, Y | \mathbf{A}, \mathbf{R}_w, \Theta) \\ = \prod_{t=1}^T f(\mathbf{x}(t) | \mathbf{s}(t), \mathbf{A}, \mathbf{R}_w) f(\mathbf{s}(t) | y(t), \Theta) \end{aligned} \quad (40)$$

Substituting (39) and (40) into (11), it is straightforward to derive that

$$\begin{aligned} J = \sum_{t=1}^T \sum_{m=1}^M \int_{\mathbf{s}} f(y(t) = m | \mathbf{s}(t), \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) \\ \log \omega_m f(\mathbf{s}(t) | y(t) = m, \Theta) d\mathbf{s} + \\ \sum_{t=1}^T \int_{\mathbf{s}} f(\mathbf{s}(t) | \mathbf{X}, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) \log f(\mathbf{x}(t) | \mathbf{s}(t), \mathbf{A}, \mathbf{R}_w) d\mathbf{s} \end{aligned}$$

APPENDIX B PROOF OF EQUATION (14)

Based on the Bayesian theory, it is easy to obtain

$$\begin{aligned} f(\mathbf{s}(t) | \mathbf{x}(t), \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) \\ = \sum_{y(t)=1}^M f(\mathbf{x}(t) | \mathbf{s}(t), \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) \\ f(\mathbf{s}(t) | y(t) = m, \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) \end{aligned} \quad (41)$$

Hence

$$\begin{aligned} f(\mathbf{x}(t), \mathbf{s}(t), \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) \\ = \frac{1}{|2\pi \mathbf{R}_w^g|^{1/2}} \\ \exp \left\{ -\frac{1}{2} (\mathbf{x}(t) - \mathbf{A}^g \mathbf{s}(t))^T (\mathbf{R}_w^g)^{-1} (\mathbf{x}(t) - \mathbf{A}^g \mathbf{s}(t)) \right\} \\ \sum_{m=1}^M \omega_m^g \frac{1}{|2\pi \mathbf{C}_m^g|^{1/2}} \\ \exp \left\{ -\frac{1}{2} (\mathbf{s}(t) - \boldsymbol{\mu}_m^g)^T (\mathbf{C}_m^g)^{-1} (\mathbf{s}(t) - \boldsymbol{\mu}_m^g) \right\} \\ = \sum_{m=1}^M \omega_m^g \frac{1}{|2\pi \mathbf{R}_w^g|^{1/2}} \frac{1}{|2\pi \mathbf{C}_m^g|^{1/2}} \\ \exp \left\{ -\frac{1}{2} (\mathbf{x}(t) - \mathbf{A}^g \mathbf{s}(t))^T (\mathbf{R}_w^g)^{-1} (\mathbf{x}(t) - \mathbf{A}^g \mathbf{s}(t)) \right\} \\ \exp \left\{ -\frac{1}{2} (\mathbf{s}(t) - \boldsymbol{\mu}_m^g)^T (\mathbf{C}_m^g)^{-1} (\mathbf{s}(t) - \boldsymbol{\mu}_m^g) \right\} \end{aligned} \quad (42)$$

After a series of derivations, one obtains

$$f(\mathbf{s}(t) | \mathbf{x}(t), \mathbf{A}^g, \mathbf{R}_w^g, \Theta^g) = \sum_{m=1}^M \tilde{\omega}_{mt}^g \mathcal{N} \left[\mathbf{s}(t); \tilde{\boldsymbol{\mu}}_{mt}^g, \tilde{\mathbf{C}}_{mt}^g \right] \quad (43)$$

where

$$\begin{cases} \tilde{\mathbf{C}}_{mt}^g = ((\mathbf{A}^g)^T (\mathbf{R}_w^g)^{-1} \mathbf{A}^g + (\mathbf{C}_m^g)^{-1})^{-1} \\ \tilde{\boldsymbol{\mu}}_{mt}^g = (\tilde{\mathbf{C}}_{mt}^g) ((\mathbf{A}^g)^T (\mathbf{R}_w^g)^{-1} \mathbf{x}(t) + (\mathbf{C}_m^g)^{-1} \boldsymbol{\mu}_m^g) \\ \tilde{\omega}_{mt}^g = \omega_m^g \left(\frac{|\tilde{\mathbf{C}}_{mt}^g|^{1/2}}{|2\pi \mathbf{R}_w^g|^{1/2} |\mathbf{C}_m^g|^{1/2}} \right) \\ \exp \left\{ -\frac{1}{2} \left[\mathbf{x}^T(t) (\mathbf{R}_w^g)^{-1} \mathbf{x}(t) \right. \right. \\ \left. \left. + (\boldsymbol{\mu}_m^g)^T (\mathbf{C}_m^g)^{-1} \boldsymbol{\mu}_m^g - (\tilde{\boldsymbol{\mu}}_{mt}^g)^T (\tilde{\mathbf{C}}_{mt}^g)^{-1} \tilde{\boldsymbol{\mu}}_{mt}^g \right] \right\} \end{cases}$$

APPENDIX C DEFINITION OF SIMILARITY SCORE

In order to measure the separation performance, the similarity score is introduced to evaluate the separation performance

of the proposed algorithm

$$\rho_{ii} = \frac{\sum_{t=1}^T s_i(t)\hat{s}_i(t)}{\sqrt{\sum_{t=1}^T (s_i(t))^2 \sum_{t=1}^T (\hat{s}_i(t))^2}} \quad (44)$$

where $\hat{s}_i(t)$ is the i th recovered source signal. ρ_{ii} depicts the similarity between the i th original source signal and the corresponding recovered source signal. It is clear that, the larger the value of ρ_{ii} , the higher the degree of similarity between the original sources and the recovered sources.

REFERENCES

- [1] A. J. Bell, and T. J. Sejnowski, "An information-maximization approach to blind separation and blind deconvolution," *Neural Computation*, vol. 7, pp. 1129-1159, 1995.
- [2] J. F. Cardoso, and B. H. Laheld, "Equivariant adaptive source separation," *IEEE Transactions on Signal Processing*, vol. 44, no. 12, pp. 3017-3030, 1996.
- [3] J. F. Cardoso, and A. Souloumiac, "Blind beamforming for non-Gaussian signals," *IEE Proceedings F on Radar and Signal Processing*, vol. 140, no. 6, pp. 362-370, 1993.
- [4] A. Hyvärinen, "Fast and robust fixed-point algorithms for independent component analysis," *IEEE Transactions on Neural Networks*, vol. 10, no. 3, pp. 626-634, 1999.
- [5] S. Amari, and A. Cichocki, "Adaptive blind signal processing-neural network approaches," *Proceedings of IEEE*, vol. 86, no. 10, pp. 2026-2048, 1998.
- [6] L. D. Lathauwer, J. Castaing, and J. F. Cardoso, "Fourth-order cumulant-based blind identification of underdetermined mixtures," *IEEE Transactions on Signal Processing*, vol. 55, no. 6, pp. 2965-2973, 2007.
- [7] L. D. Lathauwer, and J. Castaing, "Blind identification of underdetermined mixtures by simultaneous matrix diagonalization," *IEEE Transactions on Signal Processing*, vol. 56, no. 3, pp. 1096-1105, 2008.
- [8] A. Karfoul, L. Albera, and G. Birot, "Blind underdetermined mixture identification by joint canonical decomposition of HO cumulants," *IEEE Transactions on Signal Processing*, vol. 58, no. 2, pp. 638-649, 2010.
- [9] P. Comon, and M. Rajih, "Blind identification of underdetermined mixtures based on the characteristic function," *Signal Processing*, vol. 86, no. 9, pp. 2671-2681, 2006.
- [10] X. Luciani, A. L. F. de Almeida, and P. Comon, "Blind identification of underdetermined mixtures based on the characteristic function: The complex case," *IEEE Transactions on Signal Processing*, vol. 59, no. 2, pp. 540-553, 2011.
- [11] F. Gu, H. Zhang, and D. Zhu, "Blind separation of complex sources using generalised generating function," *IEEE Signal Processing Letters*, vol. 20, no. 1, pp. 71-74, 2013.
- [12] F. Gu, H. Zhang, W. Wang, and D. Zhu, "Generalized generating function with Tucker decomposition and alternating least squares for underdetermined blind identification," *EURASIP Journal on Advances in Signal Processing*, 2013, doi:10.1186/1687-6180-2013-124.
- [13] S. Peng, and W. Hwang, "Null space pursuit: An operator-based approach to adaptive signal processing," *IEEE Transactions on Signal Processing*, vol. 58, no. 5, pp. 2475-2483, 2010.
- [14] W. Hwang, J. Ho, and Y. Lin, "Null space component analysis for noisy blind source separation," Tech. Rep. TR-IIS-13-001, 2014. [online]. Available: <http://www.iis.sinica.edu.tw/page/library/TechReport/tr2013/tr13.html>.
- [15] H. Snoussi, and A. M. Djafari, "Unsupervised learning for source separation with mixture of Gaussians prior for sources and Gaussian prior for mixture coefficients," *Proceedings of the 2001 IEEE Signal Processing Society Workshop on Neural Networks for Signal Processing XI*, 2001, pp. 293-302.
- [16] K. H. Knuth, "Informed source separation: a Bayesian tutorial," in: *Proceedings of the 13th European Signal Processing Conference (EUSIPCO 2005)*, Antalya, Turkey, 2005.
- [17] Y. Zhang, X. Shi, and C. H. Chen, "A Gaussian mixture model for underdetermined independent component analysis," *Signal Processing*, vol. 86, no. 6, pp. 1538-1549, 2006.
- [18] H. Snoussi, and J. Idier, "Bayesian blind separation of generalized hyperbolic processes in noisy and underdetermined mixtures," *IEEE Transactions on Signal Processing*, vol. 54, no. 9, pp. 3257-3269, 2006.
- [19] A. Belouchrani, and J. F. Cardoso, "Maximum likelihood source separation for discrete sources," *Elsevier EUSIPCO'94*, Edinburgh, Scotland, 1994, pp. 768-771.
- [20] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *Journal of Royal Statistical Society. Series B*, vol. 39, no. 1, pp. 1-38, 1977.
- [21] J. Bilmès, "A gentle tutorial on the EM algorithm and its application to parameter estimation for Gaussian mixture and hidden Markov models," University of Berkely, Tech. Rep. ICSI-TR-97-021, 1997 [online]. Available: <http://www.citeseer.nj.nec.com/blimes98gentle.html>.
- [22] J. Q. Li, and A. R. Barron, "Mixture density estimation," in *Advances in Neural Information Processing Systems*, Cambridge, MA: MIT Press, 2000, vol. 12, pp. 279-285.
- [23] S. Kullback, and R. A. Leibler, "On information and sufficiency," *The Annals of Mathematical Statistics*, vol. 22, no. 1, pp. 79-86, 1951.
- [24] E. Moulines, J. F. Cardoso, and E. Gassiat, "Maximum likelihood for blind separation and deconvolution of noisy signals using mixture models," *International Conference on Acoustics, Speech, and Signal Processing*, Munich, Germany, 1997, pp. 3617-3620.
- [25] S. Kim, and C. D. Yoo, "Underdetermined blind source separation based on generalized Gaussian distribution," *Proceedings of the 16th IEEE Signal Processing Society Workshop on Machine Learning for Signal Processing*, Arlington, VA, 2006, pp. 103-108.
- [26] Q. Huang, J. Yang, Y. Xue, and Y. Zhou, "Temporally correlated source separation based on variational Kalman smoother," *Digital Signal Processing*, vol. 18, no. 3, pp. 422-433, 2008.
- [27] S. Sun, C. Peng, W. Hou, J. Zheng, Y. Jiang, and X. Zheng, "Blind source separation with time series variational Bayes expectation maximization algorithm," *Digital Signal Processing*, vol. 12, no.1, pp. 17-33, 2012.
- [28] K. Todros, and J. Tabrikian, "Blind separation of independent sources using Gaussian mixture model," *IEEE Transactions on Signal Processing*, vol. 55, no. 7, pp. 3645-3658, 2007.
- [29] T. Routtenberg, and J. Tabrikian, "MIMO-AR system identification and blind source separation for GMM-distributed sources," *IEEE Transactions on Signal Processing*, vol. 57, no. 5, pp. 1717-1730, 2009.
- [30] F. Gu, H. Zhang, and D. Zhu, "Blind separation of non-stationary sources using continuous density Markov models," *Digital Signal Processing*, vol. 23, no. 5, pp. 1549-1564, 2013.
- [31] T. Rydn, "EM versus chain Monte Carlo for estimation of hidden Markov models: A computational perspective," *Bayesian Analysis*, vol. 3, pp. 659-688, 2008.
- [32] J. L. Gauvain, and C. H. Lee, "Maximum a posteriori estimation for multivariate Gaussian mixture observations of chains," *IEEE Transactions on Speech and Audio Processing*, vol. 2, no. 2, pp. 291-298, 1994.
- [33] Q. Huo, and C. Chan, "Bayesian adaptive learning of the parameters of the hidden Markov model for speech recognition," HKU CSIS Tech Report TR-92-08, 1992. [Online] <http://www.csis.hku.hk/research/techreps/document/TR-92-08.pdf>
- [34] F. Gu, H. Zhang, and Y. Xiao, "A Bayesian approach to blind separation of mixed discrete sources by Gibbs sampling," *Lecture Notes on Computer Science*, vol. 6905, pp. 463-475, 2011.
- [35] [Online] <http://ai.korea.ac.kr/classes/2004/cse827/doc/map.1994.ieee.291.pdf>, accessed in 2013.
- [36] Y. Zhao, "Image segmentation using temporal-spatial information in dynamic scenes," in *Proceedings of the IEEE Int. Conf. on Machine Learning and Cybernetics*, 2003.
- [37] L. Xu, and M. I. Jordan, "On convergence properties of the EM algorithm for Gaussian mixtures," *Neural Computation*, vol. 8, pp. 129-151, 1996.
- [38] J. Ma, L. Xu, and M. I. Jordan, "Asymptotic convergence rate of the EM algorithm for Gaussian mixtures," *Neural Computation*, vol. 12, pp. 2881-2907, 2000.
- [39] P. Tseng, "Convergence of a block coordinate descent method for nondifferentiable minimization," *Journal of Optimization Theory and Applications*, vol. 109, pp. 475-494, 2001.
- [40] M. Razaviyayn, M. Hong, and Z. Luo, "A unified convergence analysis of block successive minimization methods for nonsmooth optimization," arXiv:1209.2385 [math.OA] 11 Sep 2012.
- [41] J.-L. Gauvain, "Speaker adaptation based on MAP estimation of HMM parameters," in *Proceedings of the IEEE Int. Conf. on Acoustics Speech and Signal Processing*, 1993.
- [42] G. Schwarz, "Estimation the dimension of a model," *Annals of Statistics*, vol. 6, no. 2, pp. 461-464, 1978.