



**HAL**  
open science

## EM Methods for Finite Mixtures

Gilles Celeux

► **To cite this version:**

Gilles Celeux. EM Methods for Finite Mixtures. Sylvia Frühwirth-Schnatter, Gilles Celeux, Christian P. Robert. Handbook of Mixture Analysis, CRC Press, pp.498, 2018. hal-01973069

**HAL Id: hal-01973069**

**<https://inria.hal.science/hal-01973069>**

Submitted on 8 Jan 2019

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

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

# EM Methods for Finite Mixtures

Gilles Celeux  
INRIA Saclay, France

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>The EM Algorithm</b>	<b>1</b>
2.1	Description of EM for finite mixtures . . . . .	1
2.2	EM as an alternating-maximization algorithm . . . . .	3
<b>3</b>	<b>Convergence and Behavior of EM</b>	<b>4</b>
<b>4</b>	<b>Cousin Algorithms of EM</b>	<b>5</b>
4.1	Stochastic versions of the EM algorithm . . . . .	6
4.2	The Classification EM algorithm . . . . .	8
<b>5</b>	<b>Accelerating the EM Algorithm</b>	<b>10</b>
<b>6</b>	<b>Initializing the EM Algorithm</b>	<b>11</b>
6.1	Random initialization . . . . .	11
6.2	Hierarchical initialization . . . . .	12
6.3	Recursive initialization . . . . .	12
<b>7</b>	<b>Avoiding Spurious Local Maximizers</b>	<b>12</b>
<b>8</b>	<b>Concluding Remarks</b>	<b>14</b>

## 1 Introduction

Even in the simplest situation of a two-component mixture where only the mixing proportion  $\eta_1$  is missing, the likelihood equation of a mixture model is not available in closed form. Obviously, in such a simple situation the maximum likelihood estimate (MLE) of the mixture parameters can be derived easily with standard optimization algorithms such as Newton–Raphson. But the number of parameters  $\theta$  in a mixture model grows rapidly with the dimension  $d$  of variables and with the number  $G$  of components. This means the Newton–Raphson algorithm becomes expensive both mathematically and computationally for evaluating the observed information matrix of the vector parameter  $\theta$ . Moreover, this algorithm does not increase the likelihood of being maximized at each of its iterations. It is thus no surprise that

the Newton–Raphson algorithm is far from being the most exploited algorithm to derive the MLE of a finite mixture model.

By contrast, the EM algorithm Dempster et al. (1977) stands as the most popular algorithm for the derivation of the MLE or the posterior mode for hidden structure models. Since mixture models are typical cases of hidden structure models, where the allocation variables  $z_i$ ,  $i = 1, \dots, n$ , are missing (as outlined in Chapter 1), EM applies to this setting. As a matter of fact, and as noted in Chapter 1, mixtures are often used as a prime illustration of the implementation of an EM algorithm. Indeed, mixture structures allows one to clearly highlight the rationale, the advantages, and the possible drawbacks of the EM algorithm McLachlan and Peel (2000). In this chapter, we restrict our attention to maximum likelihood estimation for mixtures, which still leads to a clear description of the EM algorithm.

## 2 The EM Algorithm

### 2.1 Description of EM for finite mixtures

Often the presence of missing data in the model of interest makes maximum likelihood (ML) inference difficult. This is clear when considering the mixture model, since the observed-data or observed log likelihood  $\mathbf{y} = (y_1, \dots, y_n)$ ,

$$\ell_o(\theta) = \sum_{i=1}^n \log \left[ \sum_{g=1}^G \eta_g f_g(y_i | \theta_g) \right],$$

involves the logarithm of a sum. The idea of the EM algorithm in this context is to make use of the missing labels  $z_i$  by maximizing at each iteration the conditional expectation of the complete-data or completed log likelihood

$$\ell_c(\theta, \mathbf{z}; G) = \sum_{i=1}^n \sum_{g=1}^G z_{ig} \log(\eta_g f_g(y_i | \theta_g)),$$

given the observed data and a current value of the parameter, towards the derivation of the MLE of its vector parameter in a simple way.

The EM algorithm takes advantage of the simple relation connecting the observed and the completed likelihoods,

$$\ell_o(\theta) = \ell_c(\theta, \mathbf{z}; G) - \sum_{g=1}^G \sum_{i=1}^n z_{ig} \log \tau_{ig},$$

where  $\tau_{ig}$  is the conditional probability that  $y_i$  arises from component  $g$ , for  $i = 1, \dots, n$  and  $g = 1, \dots, G$ , given the parameter value  $\theta$ . This leads to

$$\ell_o(\theta) = Q(\theta | \theta^{(s)}) - H(\theta | \theta^{(s)}),$$

where  $H(\theta | \theta^{(s)}) = \sum_{g=1}^G \sum_{i=1}^n \tau_{ig}^{(s)} \log \tau_{ig}$  and  $\theta^{(s)}$  is the current parameter value. From Jensen's inequality, we have that  $H(\theta | \theta^{(s)}) \leq H(\theta^{(s)} | \theta^{(s)})$ . Thus if  $Q(\theta | \theta^{(s)}) \geq Q(\theta^{(s)} | \theta^{(s)})$  then  $\ell_o(\theta) \geq \ell_o(\theta^{(s)})$  and the likelihood value increases.

Therefore, starting from an arbitrary initial value  $\theta^{(0)}$ , the EM algorithm can be summarized as follows.

**E step** Compute  $Q(\theta, \theta^{(s)}) = E(\ell_c(\theta, \mathbf{z}; G) \mid \mathbf{y}, \theta^{(s)})$ , where the expectation is taken with respect to  $p(\mathbf{z} \mid \mathbf{y}, \theta^{(s)})$ .

**M step** Determine  $\theta^{(s+1)} = \arg \max_{\theta} Q(\theta, \theta^{(s)})$ .

In the mixture context, we have

$$Q(\theta \mid \theta^{(s)}) = \sum_{i=1}^n \sum_{g=1}^G \tau_{ig}^{(s)} \{\log \eta_g + \log f_g(y_i \mid \theta_g)\},$$

$\tau_{ig}^{(s)}$  being the conditional probability that  $y_i$  arises from component  $g$ , for  $i = 1, \dots, n$  and  $g = 1, \dots, G$ , for the current parameter value  $\theta^{(s)}$ . Thus, the E step reduces to the computation of these conditional probabilities:

**E step**

$$\tau_{ig}^{(s)} = \frac{\eta_g^{(s)} f_g(y_i \mid \theta_g^{(s)})}{\sum_{g'=1}^G \eta_{g'}^{(s)} f_{g'}(y_i \mid \theta_{g'}^{(s)})}, \quad (1)$$

for  $i = 1, \dots, n$  and  $g = 1, \dots, G$ .

Note that the above M step depends on the mixture model at hand. For instance, in the case where the density  $f_g(\cdot \mid \theta_g)$  is a multivariate Gaussian density with  $\theta_g = (\mu_g, \Sigma_g)$ , where  $\mu_g$  is the mean and  $\Sigma_g$  the covariance matrix of a Gaussian distribution, this M step amounts to the updates

**M step**

$$\begin{aligned} \eta_g^{(s+1)} &= \frac{\sum_{i=1}^n \tau_{ig}^{(s)}}{n}, \\ \mu_g^{(s+1)} &= \frac{\sum_{i=1}^n \tau_{ig}^{(s)} y_i}{\sum_{i=1}^n \tau_{ig}^{(s)}}, \\ \Sigma_g^{(s+1)} &= \frac{\sum_{i=1}^n \tau_{ig}^{(s)} (y_i - \mu_g^{(s+1)})(y_i - \mu_g^{(s+1)})^\top}{\sum_{i=1}^n \tau_{ig}^{(s)}}, \end{aligned}$$

for  $g = 1, \dots, G$ .

The EM algorithm enjoys nice practical features, which explains its widespread use. The M step is often of closed form for the very reason that the EM algorithm takes the missing labels into account, in a way easy to code that does not involve numerical difficulties since it bypasses inverting large-dimensional matrices. As apparent from the formulas in the M step described above, the updated estimates are standard ML estimates where the observations are weighted with the conditional probabilities ( $\tau_{ig}$ ). For some specific mixture models, it

may still happen that the M step is not available in closed form. Examples of such situations are found in Celeux and Govaert (1995) for Gaussian mixture models where the component covariance matrices are assumed to share the same eigenvectors but have different eigenvalues. Nonetheless, the increase in the computational burden most often remains limited.

Most importantly, Dempster et al. (1977) provide a proof that the observed log likelihood is increasing at each iteration,  $\ell_o(\theta^{(s+1)}) \geq \ell_o(\theta^{(s)})$ , with equality if and only if  $Q(\theta^{(s+1)}|\theta^{(s)}) = Q(\theta^{(s)}|\theta^{(s)})$ . As remarked above, this is a direct consequence of Jensen's inequality applied to the convex function  $H(\theta | \theta^{(s)})$  and it stands as a fundamental convergence property of the EM algorithm.

## 2.2 EM as an alternating-maximization algorithm

In the mixture context, Hathaway (1986) has shown that the EM algorithm can be viewed as an alternate optimization algorithm aiming to maximize the fuzzy clustering criterion

$$\mathcal{F}_c(\theta, \mathbf{s}) = \ell_c(\theta, \mathbf{s}) + H(\mathbf{s}) \quad (2)$$

where  $\mathbf{s} = (s_{ig})_{i=1, \dots, n}^{g=1, \dots, G}$  denotes a fuzzy clustering matrix of the  $n$  observations in  $G$  clusters,

$$\ell_c(\theta, \mathbf{s}) = \sum_{i=1}^n \sum_{g=1}^G s_{ig} \log(\eta_g f_g(y_i | \theta_g))$$

is the completed log likelihood associated to  $\mathbf{s}$ , and

$$H(\mathbf{s}) = - \sum_{i=1}^n \sum_{g=1}^G s_{ig} \log s_{ig}$$

is the entropy of the fuzzy clustering matrix  $\mathbf{s}$ . More specifically:

- (a) Since  $\sum_{g=1}^G s_{ig} = 1$  for  $i = 1, \dots, n$ , maximizing  $\mathcal{F}_c(\theta, \mathbf{s})$  with respect to  $\mathbf{s}$  for fixed  $\theta$  leads, after standard Lagrangian manipulation, to

$$s_{ig} = \tau_{ig} = \frac{\eta_g f_g(y_i | \theta_g)}{\sum_{g'=1}^G \eta_{g'} f_{g'}(y_i | \theta_{g'})},$$

for  $i = 1, \dots, n$  and  $g = 1, \dots, G$ . This is exactly the E step of EM.

- (b) Maximizing  $\mathcal{F}_c(\theta, \mathbf{s})$  in  $\theta$  for a fixed value of  $\mathbf{s}$  leads to maximizing  $\ell_c(\theta, \mathbf{s})$  since  $H(\mathbf{s})$  does not depend on  $\theta$ . Thus, this is exactly the M step of EM.

This interpretation of the EM algorithm as an alternating-maximization algorithm has been extended to a general context in Neal and Hinton (1998): EM can be regarded as an alternating-optimization algorithm to maximize the criterion

$$\mathcal{F}_c(P, \theta) = E_P(\log p(\mathbf{y}, \mathbf{z} | \theta) + H(P)), \quad (3)$$

where  $P$  is a probability distribution over the space of the missing data  $\mathbf{z}$  and

$$H(P) = -E_P(\log P)$$

is the entropy of  $P$ . Moreover, the criterion  $\mathcal{F}_c$  can be related to the Kullback–Leibler divergence  $KL$  between  $P(\mathbf{z})$  and the conditional distribution  $p(\mathbf{z} | \mathbf{y}, \theta)$  of the missing data  $\mathbf{z}$  knowing the observed data  $\mathbf{y}$  and the parameter  $\theta$ :

$$\mathcal{F}_c(P, \theta) = \ell_o(\theta) + KL(P(\mathbf{z}), p(\mathbf{z} | \mathbf{y}, \theta)), \quad (4)$$

where

$$KL(P(\mathbf{z}), p(\mathbf{z} | \mathbf{y}, \theta)) = \int P(\mathbf{z}) \log \frac{P(\mathbf{z})}{p(\mathbf{z} | \mathbf{y}, \theta)} d\mathbf{z}.$$

Relation (4) can be exploited to define variational approximations of the EM algorithm when the E step is intractable (see, for instance, Govaert and Nadif (2008); or Titterington (2011)). The idea is to restrict the distribution  $p(\mathbf{z})$  to factorize with respect to well-chosen groups  $z_t$ ,  $t = 1, \dots, K$ , so that

$$P(\mathbf{z}) = \prod_{t=1}^K P(z_t) \quad (5)$$

and to seek the distribution of the form (??) for which  $KL(P(\mathbf{z}), p(\mathbf{z} | \mathbf{y}, \theta))$  is minimum. Using (3), it is easy to see that this distribution yields a lower bound on  $\mathcal{F}_c(P, \theta)$  among the distributions  $P$  of the form (??). Obviously, the factorization in (??) is chosen to ensure the tractability of this minimization.

When needed, the variational Bayesian EM algorithm is a popular algorithm to approximate the posterior distribution of parameterized models (see, for instance, Bishop (2006); Chapter 10, this volume; or Titterington (2011)).

### 3 Convergence and Behavior of EM

General results on the theoretical behavior of the EM algorithm can be found in Dempster et al. (1977) and Wu (1983). Such results cannot be detailed here, but the most important aspects are recalled in the following. From the increase of the observed log likelihood at each iteration of the EM algorithm, it follows that the MLE of  $\theta$  is a fixed point of the EM algorithm. More precisely, denoting by  $EM$  the stepwise operator of the EM algorithm (that is,  $\theta^{(s+1)} = EM(\theta^{(s)})$ ), we have the following theorem. For  $\theta^* \in \operatorname{argmax} \ell_o(\theta)$ , we have almost surely that

$$\begin{aligned} \ell_o(EM(\theta^*)) &= \ell_o(\theta^*), \\ Q(EM(\theta^*)|\theta^*) &= Q(\theta^*|\theta^*), \end{aligned}$$

and, if  $\theta^*$  is unique,  $EM(\theta^*) = \theta^*$ . Under standard regularity conditions it can further be shown that the fixed points of the EM algorithm are stationary points of the observed log likelihood. Unfortunately, these stationary points can either be local maxima or saddle-points of the observed log likelihood, meaning that the algorithm does not necessarily produce a global maximum. Obviously if  $\ell_o(\theta)$  is unimodal with a single stationary point, then  $(\theta^{(s)})$  converges towards the unique maximizer  $\theta^*$  of  $\ell_o(\theta)$ . Moreover, we have the following result. Under standard regularity conditions, any fixed point  $\theta^*$  of the  $EM$  algorithm satisfies the relation

$$D(EM)(\theta^*) = (D^{20}Q(\theta^*|\theta^*))^{-1} D^{20}H(\theta^*|\theta^*), \quad (6)$$

$D(EM)(\theta^*)$  being the Jacobian matrix of the operator  $EM$  at  $\theta^*$ , and  $D^{20}Q(\theta^*|\theta^*)$  ( $D^{20}H(\theta^*|\theta^*)$ ) being the Hessian matrix of  $Q(\theta^*|\theta^*)$  ( $H(\theta^*|\theta^*)$ ) with respect to its first argument at  $\theta^*$ .

Thus the convergence rate of the EM algorithm towards a fixed point  $\theta^*$  is determined by the eigenvalues of  $D(EM)$ . Moreover, from (5), we have

$$I(\theta|\mathbf{y}) = -D^{20}Q(\theta^*|\theta^*) + D^{20}H(\theta^*|\theta^*), \quad (7)$$

where  $I(\theta|\mathbf{y})$  is the empirical observed information on  $\theta$ ,  $-D^{20}Q(\theta^*|\theta^*)$  is complete information, and  $-D^{20}H(\theta^*|\theta^*)$  is missing information, and  $D(EM)(\theta^*)$  can be regarded as the ‘‘ratio’’ of the missing information over the complete information. The following result can also be deduced from (??). For any fixed point  $\theta^*$  of the EM algorithm,

$$D^2\ell_o(\theta^*) = D^{20}Q(\theta^*|\theta^*) [I_d - D(EM)(\theta^*)], \quad (8)$$

$I_d$  being the identity matrix. This result means that, provided the matrix  $D^{20}Q(\theta^*|\theta^*)$  is negative definite, a fixed point  $\theta^*$  of the EM algorithm is attractive (with its eigenvalues belonging to  $[0, 1]$ ) if and only if it is a local maximum of  $\ell_o$ . From (??), the greater the largest eigenvalue of  $D(EM)(\theta^*)$ , the slower the convergence of the EM algorithm towards an attractive fixed point  $\theta^*$ .

In the specific framework of mixtures with components from an exponential family, Redner and Walker (1984) refine these results to show that, under regularity conditions:

- (a) the unique solution of the log likelihood equations almost surely exists;
- (b) for  $n$  large enough ( $\theta^{(s)}$ ) linearly converges towards this solution, provided the initial position  $\theta^{(0)}$  of the EM algorithm is not too far this optimal solution.

The important regularity conditions of Redner and Walker (1984) are that the mixing proportions are positive and the Fisher information matrix when evaluated at the true  $\theta$  is positive definite.

The results of Redner and Walker (1984) are helpful in highlighting some important features of the EM algorithm, in particular concerning its practical ability to derive the ML estimates of the mixture parameters. First, as for most iterative algorithms used to optimize a non-convex function, the EM algorithm solution depends on its initial position. This dependence can be severe when the log likelihood function includes many local maxima and saddle-points, since the EM algorithm stops at the first fixed point it reaches.

Second, the rate of linear convergence to an attractive fixed point is determined by the largest eigenvalue, always smaller than 1, of  $D(EM)$ . Thus the larger the missing information, the slower the EM algorithm. In practice, the EM algorithm is well known for often having dramatically slow convergence speeds, even when the log likelihood function is convex. An example of this slow convergence is discussed in detail in Campillo and Le Gland (1989). We report in Figure 1 two different situations these authors dealt with. On the left-hand side, the missing information is low and the EM algorithm converges in five iterations. On the right-hand side, the missing information is high and the EM algorithm has not yet converged after 200 iterations, with only the first 12 iterations shown in Figure 1.

Moreover, it happens that some mixture distributions, like Gaussian distributions with free component-specific covariance matrices, have unbounded likelihood functions. In such cases, the EM algorithm is jeopardized by degenerate solutions (spurious local maximizers).

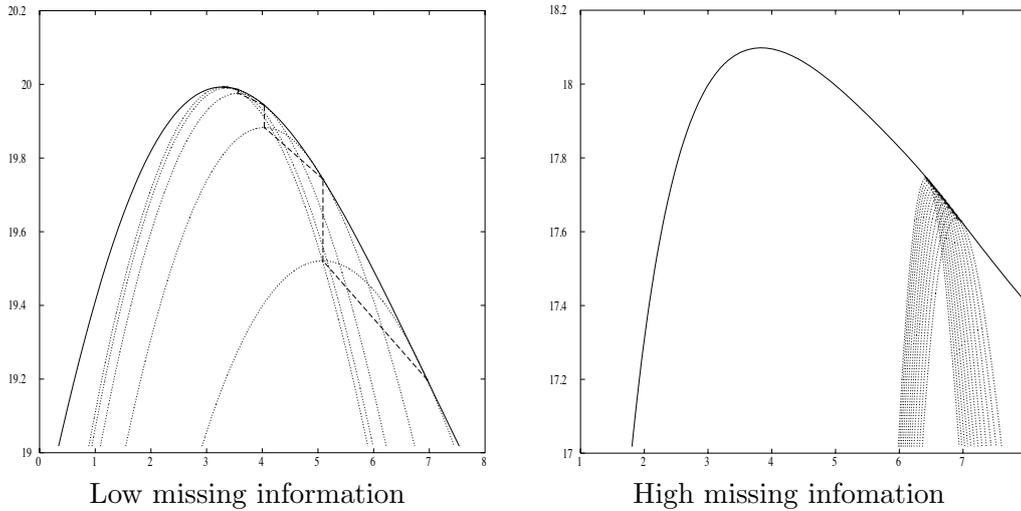


Figure 1: EM convergence behavior for a low missing information case (left) and EM convergence behavior for a high missing information case (right). The observed log likelihood function is shown as a solid line against the parameter to be estimated; the iterations of EM are shown as dashed lines, and the successive graphs of the function  $Q(\theta^{(s+1)} | \theta^{(s)})$  are shown as dotted lines.

Despite these possible drawbacks, the EM algorithm generally does a good job of deriving the MLE or the posterior mode of a mixture model. The algorithm remains, without a doubt, *the* reference algorithm for ML estimation in a mixture model. However, improved versions of EM may have to be used in order to avoid the known pitfalls of the original algorithm (slow convergence, possibly high dependence on initialization, spurious maximizers). In the following section, several ways to address these pitfalls are presented.

## 4 Cousin Algorithms of EM

All the algorithms presented in this section share with the EM algorithm the characteristic of making use of completed data by approximating the missing labels knowing the observed data and a current value of the model parameters. This is achieved in different ways, depending on the focus of interest.

### 4.1 Stochastic versions of the EM algorithm

Starting from an initial value, the Stochastic EM (SEM) algorithm Celeux and Diebolt (1985) replaces the missing labels by simulating them from their conditional distribution, knowing the observed data and a current value of the mixture parameters. Iteration  $s$  of the SEM algorithm is as follows.

**E step** The E step is the same as for EM and consists of computing the conditional probabilities  $\tau_{ig}^{(s)}$  that  $y_i$  originates from component  $g$  for  $i = 1, \dots, n$  and  $g = 1, \dots, G$ , for the

current parameter value  $\theta^{(s)}$  as in (1).

Next is the stochastic step:

**S step** Simulate the missing labels  $z_{ig}^{(s)}$  according to the multinomial distribution with parameters  $(\tau_{i1}^{(s)}, \dots, \tau_{iG}^{(s)})$ . This step results in a completed sample  $(y_i, z_i^{(s)})$ ,  $i = 1, \dots, n$ .

**M step** This step consists of maximizing the completed log likelihood  $\ell_c(\theta, \mathbf{z}; G)$  to obtain the next maximizer  $\theta^{(s+1)}$ . Typically, this step is analogous to the M step of EM. It simply leads to replacing the  $\tau_{ig}^{(s)}$  conditional probabilities with the  $z_{ig}^{(s)}$  in the likelihood equations.

Thus, SEM generates a Markov chain that is aperiodic, irreducible, and ergodic under mild conditions Diebolt and Celeux (1993). Its stationary distribution  $\Psi_n$  is approximatively centered at the ML estimator of  $\theta$ . Under regularity conditions, Nielsen (2000) showed that if  $X$  denotes a random vector drawn from the stationary distribution  $\Psi_n$  and  $\theta_0$  is the true value of the parameter  $\theta$ , then  $\sqrt{n}(X - \theta_0)$  converges to a Gaussian distribution with mean zero and regular variance matrix  $I(\theta_0)^{-1}[I_d - \{I_d - F(\theta_0)\}^{-1}]$ , where  $I(\theta_0)$  denotes the observed Fisher information matrix of  $\theta_0$ ,  $I_d$  the identity matrix, and  $F(\theta_0)$  the expected fraction of missing information.

Thus, a natural parameter estimate derived from an SEM sequence is the mean of the iterated values, typically obtained after a *burn-in* period (SEMmean). An alternative estimate is to consider the parameter value leading to the largest log likelihood in an SEM sequence (SEMmax). In practice, both pointwise estimators perform similarly.

In a mixture context, where the E step is most often simple, the main appeal of the SEM algorithm is avoiding being stuck at the first stationary point of the log likelihood function it reaches. At each iteration there is a non-zero probability of accepting an updated estimate  $\theta^{(s+1)}$  with log likelihood value lower than at  $\theta^{(s)}$ . For this very reason, the SEM algorithm avoids being stuck at a saddle-point or a spurious maximizer of the likelihood function. Thus the SEM algorithm can further be exploited to get a good starting value for the EM algorithm, if need be. Starting from SEMmean or SEMmax positions, an EM algorithm will likely converge to the MLE in a few iterations.

Different algorithms based on simulations as in the SEM algorithm are available. Some are mentioned below.

**The MCEM algorithm** This algorithm proposes a Monte Carlo implementation of the E step of the EM algorithm Wei and Tanner (1990). It replaces the computation of  $Q(\theta|\theta^{(s)})$  by the derivation of an empirical version  $Q_{(s+1)}(\theta|\theta^{(s)})$ , based on  $m$  draws of the missing vector  $\mathbf{z}$  from its current conditional distribution  $p(\mathbf{z}|\mathbf{y}, \theta^{(s)})$ . If  $m = 1$ , MCEM reduces to SEM, while for large values of  $m$ , MCEM behaves like EM. Usually, in order to achieve a compromise between speed and precision, MCEM is run in a “simulated annealing” way with small values of  $m$  in the first iterations and increasingly larger values of  $m$  as  $s$  increases. With a suitable rate of convergence of  $m$  to infinity, MCEM can be proven to converge to a sensible local maximum of the likelihood function Fort and Moulines (2003). In the mixture context, MCEM, which remains a rather slow algorithm, appears to be of little use, since, in practice,

the E step corresponds to the computation of  $Q(\theta|\theta^{(s)})$  and this reduces to the derivation of the conditional probabilities  $\tau_{ig}^{(s)}$  that  $y_i$  is generated from component  $g$ ,  $i = 1, \dots, n$  and  $g = 1, \dots, G$ , for the current parameter value  $\theta^{(s)}$ .

**The Simulated Annealing EM algorithm** This algorithm implements the intuition given above, namely, operating like an SEM algorithm at the start and finishing like an EM algorithm. More precisely, the  $s$ th iteration update of  $\theta$  with Simulated Annealing EM is

$$\theta^{(s+1)} = (1 - \gamma_{s+1})\theta_{EM}^{(s+1)} + \gamma_{s+1}\theta_{SEM}^{(s+1)}$$

where  $\theta_{EM}^{(s+1)}$  ( $\theta_{SEM}^{(s+1)}$ ) is the updated value of  $\theta$  under EM (SEM) and  $(\gamma_s)$  is a sequence of non-negative real numbers decreasing to zero with initial value  $\gamma_0 = 1$ . A slow convergence rate of  $\gamma_s$  is necessary for good performance. The conditions  $\lim_{s \rightarrow \infty} (\gamma_s / \gamma_{s+1}) = 1$  and  $\sum_s \gamma_s = \infty$  are necessary to ensure the almost sure convergence of the Simulated Annealing EM algorithm to a local maximizer of the observed log likelihood whatever the starting point. This was established in the context of finite mixtures from some exponential family by Celeux and Diebolt (1992). From a practical point of view, it is important that  $\gamma_s$  remains close to  $\gamma_0 = 1$  during the early iterations in order to allow the algorithm to avoid suboptimal stationary values of  $\ell_o(\theta)$ .

**The Stochastic Approximation EM algorithm** This algorithm is analogous to the Simulated Annealing EM algorithm and works as follows, starting from an arbitrary initial value  $\theta^{(0)}$ .

**Simulation step** As in SEM, this step makes use of  $m(s)$  realizations  $\mathbf{z}_j, j = 1, \dots, m(s)$ , of the missing label vector that are simulated from the multinomial distribution with parameters  $(\tau_{i1}^{(s)}, \dots, \tau_{iG}^{(s)})$  for  $i = 1, \dots, n$ , where  $m(s)$  is an increasing sequence of integers starting from  $m(1) = 1$ .

**Stochastic approximation step** This step computes the current approximation of  $Q(\theta, \theta^{(s)})$  according to

$$\hat{Q}_{s+1}(\theta) = \hat{Q}_s(\theta) + \gamma_s \left( \frac{1}{m(s)} \sum_{j=1}^{m(s)} \ell_c(\theta, \mathbf{z}_j) - \hat{Q}_s(\theta) \right), \quad (9)$$

where  $(\gamma_s)$  is a sequence of non-negative real numbers decreasing to zero.

**M step** This step derives  $\theta^{(s+1)} = \arg \max_{\theta} \hat{Q}_{s+1}(\theta)$ .

If the sequence  $(\gamma_s)$  is such that  $\sum_{s=1}^{\infty} \gamma_s = \infty$  and  $\sum_{s=1}^{\infty} \gamma_s^2 < \infty$ , then under regularity conditions, Delyon et al. (1999) proved that the Stochastic Approximation EM algorithm converges to a local maximum of the observed log likelihood. From a practical point of view, this algorithm is expected to behave as the Simulated Annealing EM algorithm and it is important that the sequence  $\gamma_s$  decreases slowly to zero while the number of simulations could be set to one at each iteration, that is,  $m(s) = 1$  for any  $s$ .

The goal of both the Simulated Annealing and the Stochastic Approximation EM algorithms is to provide pointwise estimates that are expected to converge to the MLE, while the SEM algorithm generates a Markov chain whose stationary distribution is expected to be centered on the MLE. In practice, there are no significant differences between the estimates provided by the Simulated Annealing or the Stochastic Approximation EM algorithms and the estimates provided by an SEM algorithm followed by a few iterations of EM.

## 4.2 The Classification EM algorithm

The algorithm now presented is particularly relevant in the model-based clustering framework where the mixture model is considered in a clustering task where each mixture component is associated to a cluster of the data; see Chapter 8 for a comprehensive review. The Classification EM (CEM) algorithm estimates both the mixture parameters and the missing labels by maximizing the complete-data log likelihood  $\ell_c(\theta, \mathbf{z}; G)$ . The  $s$ th iteration of the CEM algorithm proceeds as follows.

**E step** The E step is similar to those in the EM and SEM algorithms. It consists of computing the conditional probabilities  $\tau_{ig}^{(s)}$  that  $y_i$  arises from component  $g$  for  $i = 1, \dots, n$  and  $g = 1, \dots, G$ , given the current parameter vector  $\theta^{(s)}$ .

The classification step consists of the following derivation of labels  $\mathbf{z}^{(s)}$  at iteration  $s$ :

**C step** Assign each observation  $y_i$  to the mixture component maximizing  $\tau_{ig}^{(s)}$  over  $g = 1, \dots, G$ , that is,

$$z_{ig}^{(s)} = \begin{cases} 1, & \text{if } g = \arg \max_{g'} \tau_{ig'}^{(s)}, \\ 0, & \text{otherwise.} \end{cases}$$

**M step** This step consists of maximizing the completed log likelihood  $\ell_c(\theta, \mathbf{z}^{(s)})$  to obtain the next maximizer value  $\theta^{(s+1)}$ . This M step is formally identical to the M step of SEM.

It is important to note that the CEM algorithm does not maximize the observed-data log likelihood but the complete-data log likelihood. Thus, the CEM algorithm is expected to provide biased estimates of the mixture parameters. The more the mixture components are overlapping, the larger the bias of CEM estimates becomes. On the other hand, the CEM algorithm is a  $k$ -means type algorithm and as such converges rapidly to a fixed point. But, in most cases, it can be expected to be quite sensitive to the initial value. Moreover, the SEM algorithm can also be considered as a stochastic version of both EM and CEM algorithms, despite these algorithms not maximizing the same criterion.

Note further that Celeux and Govaert (1992) proposed a specific Simulated Annealing version of the CEM algorithm. This algorithm substitutes the E step with the AE (“annealing”) step where the conditional probabilities  $\tau_{ig}$ ,  $i = 1, \dots, n$ ;  $g = 1, \dots, G$ , are replaced by the scores

**AE step**

$$\rho_{ig}^{(s)} = \frac{\{\eta_g^{(s)} f_g(y_i | \theta_g^{(s)})\}^{1/t_s}}{\{\sum_{g'=1}^G \eta_{g'}^{(s)} f_{g'}(y_i | \theta_{g'}^{(s)})\}^{1/t_s}},$$

for  $i = 1, \dots, n$  and  $g = 1, \dots, G$ , the sequence  $(t_s)$  being a decreasing sequence of non-negative numbers starting from  $t_0 = 1$ .

For  $t_s = 1$  this algorithm is exactly the SEM algorithm and when  $t_s$  decreases from 1 to 0 as the iteration index increases, the algorithm morphs from the SEM into the CEM algorithm. In practical situations, it is recommended to have the sequence  $(t_s)$  slowly decreasing to 0. As mentioned above, the CEM algorithm and its stochastic versions are mostly useful in the model-based clustering context. However, the empirically verified fact that the CEM algorithm converges quickly to solutions that are often reasonable makes this algorithm potentially useful for maximizing the observed log likelihood. In particular, this applies to settings where the mixture components are expected to be well separated (see Celeux and Govaert (1993)) or to use the CEM algorithm to initialize the EM algorithm, as shown in Section 6.

## 5 Accelerating the EM Algorithm

Since the EM algorithm can suffer from convergence problems, numerous methods that speed up its convergence have been proposed; see (McLachlan and Krishnan 2008, Chapter 4) for an excellent review. In Chapter 3 of this volume, an acceleration procedure based on the ‘‘Aitken acceleration’’ method is presented. In this section, we only present an alternative acceleration technique specific to the mixture framework.

As expressed in (5), the rate of convergence of EM towards a stationary point of the likelihood function is governed by the fraction of missing information in the data. A large amount of missing information can induce a slow convergence of EM. Several authors have proposed variants of the EM algorithm for counteracting such slow convergence. For instance, Fessler and Hero (1995) proposed the Space-Alternating Generalized EM (SAGE) algorithm, which updates the parameters sequentially by alternating between several small hidden-data spaces defined by the algorithm designer. In the same spirit, Meng and van Dyk (1997) conceived a general Alternating Expectation-Conditional Maximization (AECM) algorithm which couples acceleration of the convergence by allowing the data augmentation scheme to vary within and between iterations with the simplification of the computation by incorporating model reduction into the M step.

We now describe the Componentwise EM algorithm for Mixtures (CEMM) algorithm, specific to the mixture context, which basically updates only one component at a time, leaving the other parameters unchanged Celeux et al. (2001). The CEMM algorithm was inspired by the SAGE algorithm. Improved convergence rates are reached by updating the parameters sequentially on small groups of observations associated with small missing data spaces rather than one large complete data space. The idea is that less informative missing data spaces lead to smaller root-convergence factors and hence faster converging algorithms. For simplicity, the CEMM algorithm is described for a multivariate Gaussian mixture model.

CEMM considers the decomposition of the parameter vector  $\theta = (\theta_g, \eta_g)_{g=1, \dots, G}$  with  $\theta_g =$

$(\mu_g, \Sigma_g)$ . Therefore the component updated at iteration  $s$  is as follows. At iteration  $(s)$ , set

$$g = s + \left\lfloor \frac{s}{G} \right\rfloor G + 1,$$

where  $\lfloor x \rfloor$  denotes the integer part of  $x$ .

**E step** Compute, for  $i = 1, \dots, n$ ,

$$\tau_{ig}^{(s)} = \frac{\eta_g^{(s)} f_g(y_i | \theta_g^{(s)})}{\sum_{g'=1}^G \eta_{g'}^{(s)} f_{g'}(y_i | \theta_{g'}^{(s)})}.$$

**M step** Set

$$\begin{aligned} \eta_g^{(s+1)} &= \frac{\sum_{i=1}^n \tau_{ig}^{(s)}}{n}, \\ \mu_g^{(s+1)} &= \frac{\sum_{i=1}^n \tau_{ig}^{(s)} y_i}{\sum_{i=1}^n \tau_{ig}^{(s)}}, \\ \Sigma_g^{(s+1)} &= \frac{\sum_{i=1}^n \tau_{ig}^{(s)} (y_i - \mu_g^{(s+1)})(y_i - \mu_g^{(s+1)})^\top}{\sum_{i=1}^n \tau_{ig}^{(s)}}, \end{aligned}$$

and, for  $\ell \neq g$ , define  $\theta_\ell^{(s+1)} = \theta_\ell^{(s)}$ .

The advantage of the CEMM algorithm over the SAGE algorithm is that it uses the new information as soon as it is available rather than waiting until all parameters have been updated. Actually, the SAGE algorithm is not exactly a componentwise algorithm because the mixing proportions are then updated at the same iteration, which involves the whole complete data structure. For this reason, it may fail to be significantly faster than the standard EM algorithm. This shows the main intuition of the componentwise EM algorithm that Celeux et al. (2001) proposed for mixtures. No iteration requires the whole complete data space as missing data space. It can therefore be expected to converge faster in various situations. However, it is important to note that with the CEMM algorithm, the mixing proportions  $\eta_g$  do not necessarily sum to 1, hence the algorithm may leave the parameter space. But it has been proven in Celeux et al. (2001) that this algorithm is guaranteed to return to the parameter space at convergence by using the proximal interpretation of EM presented in Chrétien and Hero (2008).

In Figure 2, we reproduce from Celeux et al. (2001) two numerical situations highlighting the typical behavior of the CEMM algorithm, when compared to the EM and SAGE algorithms. The graph on the left shows that the EM algorithm is faster than the CEMM algorithm when the mixture components are well separated, while the graph on the right shows that the CEMM algorithm exhibits significant improvement over the EM algorithm with overlapping mixture components. In the latter situation the EM algorithm can be expected to encounter slow convergence situations. Thus, the CEMM algorithm can be useful in some settings. An intuitive explanation is that the componentwise strategy prevents the CEMM algorithm from staying too long at critical points (typically saddle-points) where the standard EM algorithm is likely to get trapped.

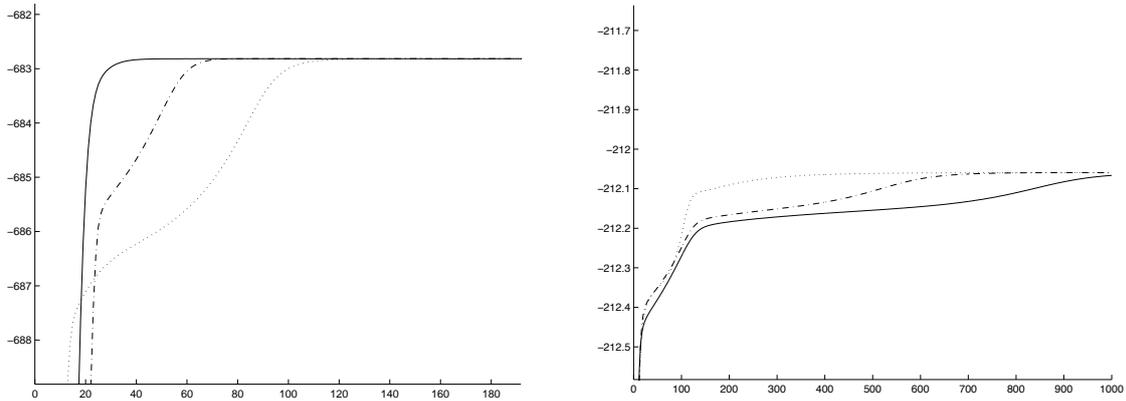


Figure 2: Comparison of the log likelihood at each cycle for EM (solid line), SAGE (dashed line) and CEMM (dotted line) for a well-separated component Gaussian mixture case (left) and in an overlapping component Gaussian mixture case (right).

## 6 Initializing the EM Algorithm

The choice of  $\theta^{(0)}$  is clearly decisive for the EM algorithm, especially when the choice of a sensible number of components  $G$  is required. In fact, information criteria used to select a mixture models are all based on the maximum likelihood values; see Chapter 7, Section 7.2.2 for a comprehensive review. Several strategies have been proposed to initialize EM for estimating the mixture parameters and they are available in most mixture software packages. Initialization strategies can be distinguished by the importance they give to randomness.

### 6.1 Random initialization

Some procedures available in the MIXMOD software (<http://www.mixmod.org>) make intensive use of random initializations and have been proposed in Biernacki et al. (2003) (see also Berchtold (2004)):

- (a) The *SEM procedure* involves starting EM with the parameter value providing the largest likelihood from a long run of the SEM algorithm.
- (b) The *CEM procedure* involves repeating the CEM algorithm a large number of times from random initial positions and starting EM with the parameter value providing the largest likelihood from these runs of the CEM algorithm.
- (c) The *small EM procedure* involves using a large number of short runs of EM. This means stopping and restarting EM at a random location before it exhibits convergence. After

a number of repetitions, EM is run from the parameter value providing the largest likelihood from these short runs of EM.

Moreover, it is recommended to try several starts from the CEM and the small EM procedure. Among these procedures, small EM is often preferred Biernacki et al. (2003), and it is the default initialization procedure in the MIXMOD software. But none of these procedures has been shown to outperform the other ones. Moreover, they can exhibit disappointing performance in high-dimensional settings because the domain parameter to be explored becomes very large. The same difficulty may appear when the number  $G$  of mixture components is large.

## 6.2 Hierarchical initialization

At the opposite end of the spectrum of initialization methods, deterministic initialization procedures such as the hierarchical procedure proposed in the `Mclust` software<sup>1</sup> do not use random starting solutions. Such hierarchical initialization can be outperformed by the above-mentioned random procedures since they do not extensively explore many initial positions. But they nonetheless provide more stable solutions and they may be preferred in high-dimensional settings or for a large number of mixture components. A recursive initialization procedure aiming at getting the best of both worlds has been proposed in Celeux and Baudry (2015) and is described in the next subsection.

## 6.3 Recursive initialization

A problem that often occurs in a mixture analysis with different numbers of components is that some solutions may be suboptimal or spurious. This problem is shared with the random initialization procedures of Section 6.1. The recursive procedures presented in Celeux and Baudry (2015) aim to avoid these irrelevant parameter estimates. Assume that the user wants to choose a mixture model with a number  $G$  of components within the range  $\{G_{\min}, \dots, G_{\max}\}$ . A recursive initialization involves splitting *at random* one of the  $G$  components into two components to get a  $(G + 1)$  solution:

1. First, the  $G_{\min}$  solution is thoroughly designed using, for instance, the *small EM* procedure repeated a large number of times.
2. From  $G = G_{\min}$ , the initial position of the  $(G + 1)$ -component mixture results from splitting one of the  $G$ -component mixture into two components.

The resulting procedures differ in the way the component to be split is chosen. In Celeux and Baudry (2015), three such strategies are considered.

- (a) *Random choice* (RC): pick the split mixture component at random Papastamoulis et al. (2016).
- (b) *Optimal sequential choice* (OSC): pick the split mixture component by optimizing a splitting criterion.

---

<sup>1</sup><http://www.stat.washington.edu/mclust/>

- (c) *Complete choice* (CC): split all  $G$  components into two components and pick the split mixture component leading to the largest likelihood.

Other splitting strategies, such as the one in Fraley et al. (2005) which is designed to deal with large data sets, are possible. Concerning strategy OSC, many sensible splitting criteria are possible, but numerical experiments reported in Celeux and Baudry (2015) show that there is little incentive to choose among these criteria. Moreover, the same numerical experiments show that strategy CC outperforms the other ones, while strategy OSC shows no advantage compared to strategy RC. Finally, CC is not so expensive and can thus be recommended.

## 7 Avoiding Spurious Local Maximizers

Deriving the maximum likelihood parameter estimate of a Gaussian mixture faces an important difficulty since the likelihood function of Gaussian mixtures with unrestricted component covariance matrices is unbounded. Thus, spurious local maximizers of the likelihood are found in the parameter space. The EM algorithm can thus fail because of these singularities depending on the starting values, models, and numbers of components, McLachlan and Peel (2000) Section 3.10, for instance. Such spurious maximizers can be avoided by imposing some constraints on the mixture parameters. The simplest and most popular constraints assume equal mixing proportions or component covariance matrices. But such constraints may sound unrealistic and cannot be regarded as a general answer to the spurious maximizer problem.

Bayesian regularization can be seen as a more general solution of this issue; see Ciuperca et al. (2003) or Fraley and Raftery (2007). It involves replacing the MLE with the maximum *a posteriori* (MAP) estimate, obtained by maximizing the penalized log likelihood  $\ell_o(\theta) + \log p(\theta)$ , with  $p(\theta)$  being a prior distribution on the vector parameter  $\theta$ . Since the Bayesian framework is only introduced to avoid degeneracies in the estimation of the component covariance matrices, there is no need to design prior distributions on the mixing proportions and the component mean vectors.

Following Fraley and Raftery (2007), it is convenient to use an inverse Wishart exchangeable conjugate prior distribution for the component-specific covariance matrices, that is,  $\Sigma_g \sim \mathcal{W}^{-1}(\nu, \Lambda)$  for  $g = 1, \dots, G$ . Obviously, the choice of the hyperparameters  $\nu$  and  $\Lambda$  is important. The choice  $\nu = d + 2$  was recommended in Fraley and Raftery (2007) and applied in Celeux and Baudry (2015). In Fraley and Raftery (2007), the scale matrix  $\Lambda$  was set to

$$\Lambda = \frac{1}{G^{1/d}} S_y, \quad (10)$$

with  $S_y$  being the empirical covariance matrix of the data  $\mathbf{y}$ . Another choice, advocated in Celeux and Baudry (2015), is

$$\Lambda = \frac{\sigma_0^{1/d}}{|S_y|^{1/d}} S_y,$$

where  $\sigma_0$  is a small positive number. This choice implies that  $|\Lambda| = \sigma_0$ . The larger  $\sigma_0$ , the larger the regularization becomes. Thus, this tuning parameter allows for control of the regularization, and weaker regularization than the fixed hyperparameter  $\Lambda$  defined in (7).

With such choices, the MAP estimate of  $\theta$  is derived via the EM algorithm. In this framework, the formulas of the EM algorithm are indeed unchanged, except for the updating of the

covariance matrices in the M step. It then becomes (at the  $(s + 1)$ th iteration):

$$\Sigma_g^{(s+1)} = \frac{\Lambda + \sum_{i=1}^n \tau_{ig}^{(s)} (y_i - \mu_g^{(s+1)})(y_i - \mu_g^{(s+1)})^\top}{\nu + \sum_{i=1}^n \tau_{ig}^{(s)} + d + 2}$$

Fraley and Raftery (2007) for details. Assuming diagonal or spherical component covariance matrices does not guard against spurious maximizers, and regularization of the covariance matrices is also desirable in these settings. For diagonal component covariance matrices  $\Sigma_g = B_{g1}, \dots, B_{gd}$ , for instance, the conjugate inverse gamma prior  $B_{gj} \sim \mathcal{IG}(\nu/2, \zeta_j/2)$  can be used for the diagonal elements  $B_{gj}$ ,  $j = 1, \dots, d$ ,  $g = 1, \dots, G$ , with hyperparameters  $\nu = d + 2$  and

$$\zeta_j = (\sigma_0)^{1/d} \frac{s_j}{s_1 \cdots s_d},$$

where  $s_j$  is the empirical variance of the variable  $j$  and  $\sigma_0$  is a small positive number. The updating of the  $B_{gj}$ s at the  $(s + 1)$ th iteration of the EM algorithm is as follows:

$$B_{gj}^{(s+1)} = \frac{\zeta_j + \sum_{i=1}^n \tau_{ig}^{(s)} (y_{ij} - \mu_{gj}^{(s+1)})^2}{\nu + \sum_{i=1}^n \tau_{ig}^{(s)} + 2}.$$

Finally, with covariance matrices  $\Sigma_g = \lambda_g I$  being proportional to the identity matrix in each component  $g = 1, \dots, G$ , the conjugate inverse gamma prior  $\lambda_g \sim \mathcal{IG}(\nu/2, \zeta/2)$  can be applied, with the hyperparameters  $\nu = d + 2$  and  $\zeta = 2(\sigma_0)^{1/d}$ , with  $\sigma_0$  again being a small positive number. The updating of the  $\lambda_g$  at the  $(s + 1)$ th iteration of the EM algorithm is as follows:

$$\lambda_g^{(s+1)} = \frac{\zeta + \sum_{i=1}^n \tau_{ig}^{(s)} (y_i - \mu_g^{(s+1)})^\top (y_i - \mu_g^{(s+1)})}{\nu + d \sum_{i=1}^n \tau_{ig}^{(s)} + d + 2}.$$

Obviously, in each case, the choice of the hyperparameter  $\sigma_0$  ends up being quite influential. For instance, it is important that the  $\sigma_0$  thus chosen does not hide the data structure. This hyperparameter is indeed driving the strength of the regularization. A careful sensitivity analysis must then be performed to choose values  $\sigma_0$  ensuring stable and meaningful estimates.

## 8 Concluding Remarks

The EM algorithm is a well-established algorithm and indeed *the* algorithm of reference for deriving the MLE or performing posterior mode estimation of mixture models. We expect it to retain its dominating position for many years to come. Indeed, Chapter 3 presents an expansion of the EM algorithm to a more general mathematical structure (the *product-of-sum* structure). Despite many of its characteristics being well documented, research on the EM and related algorithms remains very active. An important question which is still open concerns deciding when a fixed point of the EM is near the global optimum or a consistent solution of the likelihood equations, or a poor local optimum of the likelihood. With this question in mind, Balakrishnam et al. (2017) developed a theoretical framework for quantifying when and how quickly the EM algorithm converges to a small neighborhood of a given global optimum of the population likelihood (obtained in the limit with infinite data). In particular, they give precise characterizations of the region of convergence for symmetric mixtures of two Gaussians with isotropic covariance matrices. This important theoretical work confirms that to guarantee a well-behaved EM algorithm, we require a good initialization.

## References

- Balakrishnam, S., Wainwright, M. J., and Yu, B. (2017). Statistical guarantees for the EM algorithm: From population to sample-based analysis. *Annals of Statistics*, 45:77–120.
- Berchtold, A. (2004). Optimisation of mixture models: Comparison of different strategies. *Computational Statistics*, 19:385–406.
- Biernacki, C., Celeux, G., and Govaert, G. (2003). Choosing starting values for the EM algorithm for getting the highest likelihood in multivariate Gaussian mixture models. *Computational Statistics and Data Analysis*, 41:561–575.
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer-Verlag, New York.
- Campillo, F. and Le Gland, F. (1989). MLE for partially observed diffusions: Direct maximization vs. the EM algorithm. *Stochastic Processes and Applications*, 33:245–274.
- Celeux, G. and Baudry, J.-P. (2015). EM for mixtures – random initialization could be hazardous. *Statistics and Computing*, 25:713–726.
- Celeux, G., Chrétien, S., Forbes, F., and Mkhadri, A. (2001). A component-wise EM algorithm for mixtures. *Journal of Computational and Graphical Statistics*, 10:699–712.
- Celeux, G. and Diebolt, J. (1985). The SEM algorithm: a probabilistic teacher algorithm derived from the EM algorithm for the mixture problem. *Computational Statistics Quarterly*, 2:73–82.
- Celeux, G. and Diebolt, J. (1992). A Stochastic Approximation type EM algorithm for the mixture problem. *Stochastics and Stochastics Reports*, 41:119–134.
- Celeux, G. and Govaert, G. (1992). A classification EM algorithm for clustering and two stochastic versions. *Computational Statistics and Data Analysis*, 14:315–332.
- Celeux, G. and Govaert, G. (1993). Comparison of the mixture and the classification maximum likelihood in cluster analysis. *Journal of Statistical Computation and Simulation*, 47:127–146.
- Celeux, G. and Govaert, G. (1995). Gaussian parsimonious clustering models. *Pattern Recognition*, 28:781–793.
- Chrétien, S. and Hero, A. O. (2008). On EM algorithms and their proximal generalizations. *ESAIM: Probability and Statistics*, 12:308–326.
- Ciuperca, G., Ridolfi, A., and Idier, J. (2003). Penalized maximum likelihood estimator for normal mixtures. *Scandinavian Journal of Statistics*, 30:45–59.
- Delyon, B., Lavielle, M., and Moulines, E. (1999). On a stochastic approximation version of the EM algorithm. *Annals of Statistics*, 27:94–128.
- Dempster, A., Laird, N., and Rubin, D. (1977). Maximum likelihood from incomplete data via the EM algorithm (with discussion). *Journal of the Royal Statistical Society Series B*, 39:1–38.

- Diebolt, J. and Celeux, G. (1993). Asymptotic properties of a stochastic EM algorithm for estimating mixing proportions. *Communications in Statistics: Stochastic Models*, 9:599–613.
- Fessler, J. A. and Hero, A. O. (1995). Penalized maximum-likelihood image reconstruction using space-alternating generalized EM algorithms. *IEEE Transactions in Image Processing*, 4(10):1417–29.
- Fort, G. and Moulines, E. (2003). Convergence of the Monte Carlo EM for curved exponential families. *Annals of Statistics*, 31:1220–1259.
- Fraley, C., Raftery, A., and Wehrens, R. J. (2005). Incremental model-based clustering for large datasets with small clusters. *Journal of Computational and Graphical Statistics*, 14:529–546.
- Fraley, C. and Raftery, A. E. (2007). Bayesian regularization for normal mixture estimation and model-based clustering. *Journal of Classification*, 24:155–181.
- Govaert, G. and Nadif, M. (2008). Block clustering with Bernoulli mixture models: Comparison of different approaches. *Computational Statistics and Data Analysis*, 52:3233–3245.
- Hathaway, R. J. (1986). Another interpretation of the EM algorithm for mixture distributions. *Statistics and Probability Letters*, 4:53–56.
- McLachlan, G. and Peel, D. (2000). *Finite Mixture Models*. John Wiley, New York.
- McLachlan, G. J. and Krishnan, T. (2008). *The EM Algorithm and Extensions*. John Wiley, New York. Second Edition.
- Meng, X.-L. and van Dyk, D. A. (1997). The EM algorithm – an old folk-song sung to a fast new tune (with discussion). *Journal of the Royal Statistical Society Series B*, 59:511–567.
- Neal, R. M. and Hinton, G. E. (1998). A view of the EM algorithm that justifies incremental, sparse and other variants. In Jordan, M. I., editor, *Learning in Graphical Models*, pages 355–358. Kluwer Academic Publishers, Dordrecht.
- Nielsen, S. F. (2000). The stochastic EM algorithm: Estimation and asymptotic results. *Bernoulli*, 6(3):457–489.
- Papastamoulis, P., Martin-Magniette, M.-L., and Maugis-Rabusseau, C. (2016). On the estimation of mixtures of Poisson regression models with large numbers of components. *Computational Statistics and Data Analysis*, 93:97–106.
- Redner, R. and Walker, H. (1984). Mixture densities, maximum likelihood and the EM algorithm. *SIAM Reviews*, 26:195–239.
- Titterton, D. M. (2011). The EM algorithm, variational approximations and expectation propagation for mixtures. In Mengersen, K. L., Robert, C. P., and Titterton, D. M., editors, *Mixtures: Estimation and Applications*, pages 1–29. Wiley, Chichester.
- Wei, G. and Tanner, M. (1990). A Monte Carlo implementation of the EM algorithm and the poor man’s data augmentation algorithm. *Journal of the American Statistical Association*, 85:699–704.

Wu, C. (1983). On the convergence properties of the EM algorithm. *Annals of Statistics*, 11:95–103.