

## Asymptotic Convergence Rate of the EM Algorithm for Gaussian Mixtures

**Jinwen Ma**

Department of Computer Science & Engineering, The Chinese University of Hong Kong, Shatin Hong Kong and Institute of Mathematics, Shantou University, Shantou, Guangdong, 515063, People's Republic of China

**Lei Xu**

Department of Computer Science & Engineering, The Chinese University of Hong Kong, Shatin Hong Kong, People's Republic of China

**Michael I. Jordan**

Department of Computer Science and Department of Statistics, University of California at Berkeley, Berkeley, CA 94720, U.S.A.

It is well known that the convergence rate of the expectation-maximization (EM) algorithm can be faster than those of conventional first-order iterative algorithms when the overlap in the given mixture is small. But this argument has not been mathematically proved yet. This article studies this problem asymptotically in the setting of gaussian mixtures under the theoretical framework of Xu and Jordan (1996). It has been proved that the asymptotic convergence rate of the EM algorithm for gaussian mixtures locally around the true solution  $\theta^*$  is  $O(e^{-\epsilon \theta^*})$ , where  $\epsilon > 0$  is an arbitrarily small number,  $O(\cdot)$  means that it is a higher-order infinitesimal as  $\epsilon \rightarrow 0$ , and  $e^{-\epsilon \theta^*}$  is a measure of the average overlap of gaussians in the mixture. In other words, the large sample local convergence rate for the EM algorithm tends to be asymptotically superlinear when  $e^{-\epsilon \theta^*}$  tends to zero.

### 1 Introduction ---

The expectation-maximization (EM) algorithm is a general methodology for maximum likelihood (ML) or maximum a posteriori (MAP) estimation (Dempster, Laird, & Rubin, 1977). A substantial literature has been devoted to the study of the convergence of EM and related methods (e.g., Wu, 1983; Redner & Walker, 1984; Meng, 1994; Liu & Rubin, 1994; Lange, 1995a; Meng & van Dyk, 1997; Delyon, Lavielle, & Moulines, 1999). The starting point for many of these studies is the fact that EM is generally a first-order or linearly convergent algorithm, as can readily be seen by considering EM as a mapping  $\theta^{k+1} = M(\theta^k)$ , with fixed point  $\theta^* = M(\theta^*)$ . Results on

the rate of convergence of EM are obtained by calculating the information matrices for the missing data and the observed data (Dempster et al., 1977; Meng, 1994).

The theoretical convergence results obtained to date are of satisfying generality but of limited value in understanding why EM may converge slowly or rapidly in a particular problem. The existing methods to enhance the convergence rate of EM are generally based on the conventional superlinear optimization theory (e.g., Lange, 1995b; Jamshidian & Jennrich, 1997; Meng & van Dyk, 1998), and are usually rather more complex than EM. However, they are prey to other disadvantages, including an awkwardness at handling constraints on parameters. More fundamentally, it is not clear what the underlying factors are that slow EM's convergence.

It is also worth noting that EM has been successfully applied to large-scale problems such as hidden Markov models (Rabiner, 1989), probabilistic decision trees (Jordan & Jacobs, 1994), and mixtures of experts architectures (Jordan & Xu, 1995), where its empirical convergence rate can be significantly faster than those of conventional first-order iterative algorithms (i.e., gradient ascent). These empirical studies show that EM can be slow if the overlap in the given mixture is large but rapid if the overlap in the given mixture is small.

A recent analysis by Xu and Jordan (1996) provides some insight into the convergence rate of EM in the setting of gaussian mixtures. For the convenience of mathematical analyses, they studied a variant of the original EM algorithm for gaussian mixtures and showed that the condition number associated with this variant EM algorithm is guaranteed to be smaller than the condition number associated with gradient ascent, providing a general guarantee of the dominance of this variant EM algorithm over the gradient algorithm. Moreover, in cases in which the mixture components are well separated, they showed that the condition number for this EM algorithm approximately converges to one, corresponding to a local superlinear convergence rate. Thus, in this restrictive case, this type of EM algorithm has the favorable property of showing quasi-Newton behavior as it nears the ML or MAP solution. Xu (1997) further showed that the original EM algorithm has the same convergence properties as this variant EM algorithm.

We have further found by experiments that the convergence rate of the EM algorithm for gaussian mixtures is dominated by a measure of the average overlap of gaussians in the mixture. Actually, when the average overlap measure becomes small, the EM algorithm converges quickly when it nears the true solution. As the measure of the average overlap tends to zero, the EM algorithm tends to demonstrate a quasi-Newton behavior.

In this article, based on the mathematical connection between the variant EM algorithm and gradient algorithms and on one of its intermediate result on the convergence rate by Xu and Jordan (1996), as well as on the same convergence result of the original EM algorithm (Xu, 1997), we present a further theoretical analysis of the asymptotic convergence rate of the EM algorithm

locally around the true solution  $\theta^*$  with respect to  $\epsilon$ .  $\theta^*$ , a measure of the average overlap of the gaussians in the mixture. We prove a theorem that shows the asymptotic convergence rate is a higher-order in  $\epsilon$  than  $\epsilon^{0.5}$ .  $\theta^*$  when  $\epsilon$ .  $\theta^*$  tends to zero under certain conditions, where  $\epsilon > 0$  is an arbitrarily small number. Thus, we see that the large sample local convergence rate for the EM algorithm tends to be asymptotically superlinear when  $\epsilon$ .  $\theta^*$  tends to zero.

Section 2 presents the EM algorithm and the Hessian matrix for the gaussian mixture model. In section 3, we intuitively introduce our theorem on the asymptotic convergence rate of EM. Section 4 describes several lemmas needed for the proof; the proof is contained in section 5. Section 6 presents our conclusions.

## 2 The EM Algorithm and the Hessian Matrix of the Log-Likelihood

We consider the following gaussian mixture model:

$$P(x|\theta) = \sum_{j=1}^K \pi_j P(x|m_j; \Sigma_j); \quad \pi_j \geq 0; \sum_{j=1}^K \pi_j = 1; \tag{2.1}$$

where

$$P(x|m_j; \Sigma_j) = \frac{1}{(2\pi)^{d/2} |\Sigma_j|^{1/2}} e^{-\frac{1}{2}(x-m_j)^T \Sigma_j^{-1} (x-m_j)} \tag{2.2}$$

and where  $K$  is the number of the mixture components,  $x$  denotes a sample vector, and  $d$  is the dimensionality of  $x$ . The parameter vector  $\theta$  consists of the mixing proportions  $\pi_j$ , the mean vectors  $m_j$ , and the covariance matrices  $\Sigma_j = \text{cov}(x|j)$ , which are assumed positive definite.

Given  $K$  and given independently and identically distributed (i.i.d.) samples  $\{x^t\}_{t=1}^N$ , we estimate  $\theta$  by maximizing the log-likelihood:

$$l(\theta) = \log \prod_{t=1}^N P(x^t|\theta) = \sum_{t=1}^N \log P(x^t|\theta); \tag{2.3}$$

This log-likelihood can be optimized iteratively via the EM algorithm as follows:

$$\pi_j^{k+1} = \frac{1}{N} \sum_{t=1}^N h_j^{k,t} \tag{2.4}$$

$$m_j^{k+1} = \frac{1}{\sum_{t=1}^N h_j^{k,t}} \sum_{t=1}^N h_j^{k,t} \cdot x^t \tag{2.5}$$

$$\hat{\phi}_j^{k+1/} = \frac{1}{\sum_{t=1}^N h_j^{k/} \cdot t/} \sum_{t=1}^N h_j^{k/} \cdot t/ \cdot x^{t/} - m_j^{k+1/} / \cdot x - m_j^{k+1/} / T; \tag{2.6}$$

where the posterior probabilities  $h_j^{k/} \cdot t/$  are given by

$$h_j^{k/} \cdot t/ = \frac{\otimes_j^{k/} P \cdot x^{t/} | m_j^{k/}; \hat{\phi}_j^{k/}}{\sum_{i=1}^K \otimes_i^{k/} P \cdot x^{t/} | m_i^{k/}; \hat{\phi}_i^{k/}}; \tag{2.7}$$

This iterative procedure converges to a local maximum of the log-likelihood (Dempster, et al., 1977). We suppose that  $\hat{\Delta}$  is a local solution to the likelihood equation 2.3, and the EM algorithm converges to it. We now analyze the local convergence rate around this solution.

For the convenience of mathematical analyses, Xu and Jordan (1996) studied a variant of the EM algorithm by letting equation 2.6 be replaced by

$$\hat{\phi}_j^{k+1/} = \frac{1}{\sum_{t=1}^N h_j^{k/} \cdot t/} \sum_{t=1}^N h_j^{k/} \cdot t/ \cdot x^{t/} - m_j^{k/} / \cdot x - m_j^{k/} / T; \tag{2.8}$$

that is, the update of  $\hat{\phi}_j^{k+1/}$  is based on the last update  $m_j^{k/}$  instead of  $m_j^{k+1/}$  in equation 2.6. For clarity, we denote this variant of EM by VEM in this article. They showed that at each iteration, the following relationship holds between the gradient of the log-likelihood and the VEM update step:

$$\mathcal{A}^{k+1/} - \mathcal{A}^{k/} = P_{\mathcal{A}^{k/}} \frac{\partial l}{\partial \mathcal{A}} |_{\mathcal{A}=\mathcal{A}^{k/}} \tag{2.9}$$

$$m_j^{k+1/} - m_j^{k/} = P_{m_j^{k/}} \frac{\partial l}{\partial m_j} |_{m_j=m_j^{k/}} \tag{2.10}$$

$$\text{vec} [\hat{\phi}_j^{k+1/}] - \text{vec} [\hat{\phi}_j^{k/}] = P_{\hat{\phi}_j^{k/}} \frac{\partial l}{\partial \text{vec}[\hat{\phi}_j]} |_{\hat{\phi}_j=\hat{\phi}_j^{k/}}; \tag{2.11}$$

where

$$P_{\mathcal{A}^{k/}} = \frac{1}{N} \left( \text{diag} [\otimes_1^{k/}; \dots; \otimes_K^{k/}] - \mathcal{A}^{k/} \mathcal{A}^{k/T} \right) \tag{2.12}$$

$$P_{m_j^{k/}} = \frac{1}{\sum_{t=1}^N h_j^{k/} \cdot t/} \hat{\phi}_j^{k/} \tag{2.13}$$

$$P_{\hat{\phi}_j^{k/}} = \frac{2}{\sum_{t=1}^N h_j^{k/} \cdot t/} \left( \hat{\phi}_j^{k/} \otimes \hat{\phi}_j^{k/} \right); \tag{2.14}$$

and where  $\mathcal{A}$  denotes the vector of mixing proportions  $[\otimes_1; \dots; \otimes_K]^T$ ,  $j$  indexes the mixture components  $.j = 1; \dots; K/$ ,  $k$  denotes the iteration number,  $\text{vec}[B]$  denotes the vector obtained by stacking the column vectors of

the matrix  $B$ ,  $\text{vec}[B]^T = \text{vec}[B]^T$ , and  $\otimes$  denotes the Kronecker product. Moreover, given the constraints  $\sum_{j=1}^K \pi_j^{(k)} = 1$  and  $\pi_j^{(k)} \geq 0$ ;  $P_{\mathcal{A}}^{(k)}$  is a positive definite matrix and the matrices  $P_{m_j}^{(k)}$  and  $P_{\delta_j}^{(k)}$  are positive definite with probability one for  $N$  sufficiently large. Assembling these matrices into a single matrix  $P_{\mathcal{Z}} = \text{diag}[P_{\mathcal{A}}; P_{m_1}; \dots; P_{m_K}; P_{\delta_1}; \dots; P_{\delta_K}]$ ; we have the following equation for the VEM update step given by equation 3.7 in Xu and Jordan (1996):

$$\mathcal{Z}^{(k+1)} = \mathcal{Z}^{(k)} + P(\mathcal{Z}^{(k)}) \frac{\partial \ell}{\partial \mathcal{Z}} \Big|_{\mathcal{Z}=\mathcal{Z}^{(k)}} \tag{2.15}$$

where  $\mathcal{Z}$  is the collection of mixture parameters:

$$\mathcal{Z} = [\mathcal{A}^T; m_1^T; \dots; m_K^T; \text{vec}[\delta_1]^T; \dots; \text{vec}[\delta_K]^T]^T$$

In order to represent  $\mathcal{Z}$  to a set of independent variables for derivation, we introduce the following subspace,

$$\mathcal{R}_1 = \left\{ \mathcal{Z}: \sum_{j=1}^K \pi_j = 0; \pi_{pq}^{j'} = \pi_{qp}^{j'}; \text{ for all } j; p; q \right\};$$

which is obtained from

$$\mathcal{R}_2 = \left\{ \mathcal{Z}: \sum_{j=1}^K \pi_j = 1; \pi_{pq}^{j'} = \pi_{qp}^{j'}; \text{ for all } j; p; q \right\}$$

by the constant shift  $\mathcal{Z}_0$ . For the gaussian mixture, the constraint that all  $\delta_j$  are positive definite should also be added to  $\mathcal{R}_2$  and thus  $\mathcal{R}_1$ . It can be easily verified that this constraint makes  $\mathcal{R}_1$  be an open convex set of it. Since we will consider only the local differential properties of log-likelihood function at an interior point of the open convex set, we can set a new coordinate system for the parameter vector  $\mathcal{Z}$  via a set of the unit basis vectors  $E = [e_1; \dots; e_m]$ , where  $m$  is the dimension of  $\mathcal{R}_1$ .

In fact, for each  $\mathcal{Z}' = \mathcal{Z} - \mathcal{Z}_0 \in \mathcal{R}_1$ , let its coordinates under the bases  $e_1; \dots; e_m$  be denoted by  $\mathcal{Z}_c$ ; we have

$$\mathcal{Z} - \mathcal{Z}_0 = E\mathcal{Z}_c; \text{ or } \mathcal{Z} = E\mathcal{Z}_c + \mathcal{Z}_0 \tag{2.16}$$

Multiplying its both sides by  $E^T$ , it follows from  $E^T E = I$  that

$$E^T \mathcal{Z} = E^T E\mathcal{Z}_c + E^T \mathcal{Z}_0; \text{ or } \mathcal{Z}_c = E^T \mathcal{Z} - E^T \mathcal{Z}_0 \tag{2.17}$$

Putting it into equation 2.16, we have

$$\mathcal{Z}' = \mathcal{Z} - \mathcal{Z}_0 = E^T E \cdot \mathcal{Z} - \mathcal{Z}_0 = E E^T \mathcal{Z}' \text{ for } \mathcal{Z}' \in \mathcal{R}_1; \tag{2.18}$$

Thus for a matrix  $A$ , let  $\|A\| = \max_{\|2'\|=1; 2' \in \mathcal{R}_1} \|A2'\|$  be its norm constrained on  $\mathcal{R}_1$ . From equation 2.18, we certainly have the equality that  $\|A\| = \|AEE^T\|$ . We use the Euclidean norm for vectors in this article and have that

$$\begin{aligned} \|E^T\| &= \max_{\|2\|=1; 2 \in \mathcal{R}_1} \|E^T 2\| = \max_{\|2\|=1; 2 \in \mathcal{R}_1} .2^T E E^T 2 / \sqrt{2} \\ &= \max_{\|2\|=1; 2 \in \mathcal{R}_1} .2^T 2 / \sqrt{2} = \max_{\|2\|=1; 2 \in \mathcal{R}_1} \|2\| = 1: \end{aligned}$$

Following the last inequality on p. 137 of Xu and Jordan (1996), we have the local convergence rate around  $\hat{2}$  is bounded by

$$\begin{aligned} r &= \lim_{k \rightarrow \infty} \frac{\|2^{.k+1/} - \hat{2}\|}{\|2^{.k/} - \hat{2}\|} \\ &\leq \|E^T . I + P . \hat{2} / H . \hat{2} // \| = \|E^T . I + P . \hat{2} / H . \hat{2} // E E^T \| \\ &\leq \|E^T . I + P . \hat{2} / H . \hat{2} // E \| \|E^T \|: \end{aligned} \tag{2.19}$$

By the fact  $\|E^T\| = 1$ , we have

$$r \leq \|I + E^T P . \hat{2} / H . \hat{2} / E \|: \tag{2.20}$$

For the original EM algorithm, Xu (1997) further showed that the convergence rate by the original EM algorithm and the VEM algorithm is the same. Therefore, equation 2.20 also holds for the original EM algorithm. As a result, the following analyses and results apply to both the EM and VEM algorithm.

Suppose that the samples  $\{x^{.t/}\}_1^N$  are randomly selected from the gaussian mixture with the parameter  $2^*$ . When the gaussian mixture model satisfies certain regularity conditions, the EM iterations arrive at a consistent solution on maximizing log-likelihood equation 2.3 (Veaux, 1986; Redner & Walker, 1984). In this article, we assume that the EM algorithm asymptotically converges to this true parameter correctly (i.e., when  $N$  is large, for the sample data  $\{x^{.t/}\}_1^N$ , the EM algorithm converges to  $\hat{2}$  with  $\lim_{N \rightarrow \infty} \hat{2} = 2^*$ ), and we analyze the local convergence rate around this consistent solution in the limit form. It follows from equation 2.20 that an upper bound of the asymptotic convergence rate is given by

$$\begin{aligned} r &\leq \lim_{N \rightarrow \infty} \|I + E^T P . \hat{2} / H . \hat{2} / E \| \\ &= \|I + E^T \lim_{N \rightarrow \infty} P . \hat{2} / H . \hat{2} / E \| \\ &= \|I + E^T \lim_{N \rightarrow \infty} P . 2^* / H . 2^* / E \|: \end{aligned} \tag{2.21}$$

In order to estimate this bound, the rest of this article pursues an analysis of the convergence behavior of  $P_{.2^*}/H_{.2^*}$  as  $N$  increases to infinity.

First, we give the Hessian of the log-likelihood equation 2.3 in the block forms. The detailed derivation is omitted (for an example of such a derivation, see Xu and Jordan, 1996, and for the general methods for matrix derivatives, see Gerald, 1980):

$$H_{\mathcal{A};\mathcal{A}^T} \triangleq \frac{\partial^2 l}{\partial \mathcal{A} \partial \mathcal{A}^T} = - \sum_{t=1}^N H_A^{-1} \cdot t / H_A^{-1} \cdot t //^T;$$

$$H_{\mathcal{A};m_j^T} \triangleq \frac{\partial^2 l}{\partial \mathcal{A} \partial m_j^T} = \sum_{t=1}^N R_{A_j}^{-1} \cdot t / x \cdot t' - m_j /^T \delta_j^{-1};$$

$$H_{\mathcal{A};\delta_j^T} \triangleq \frac{\partial^2 l}{\partial \mathcal{A} \partial \delta_j^T} = -\frac{1}{2} \sum_{t=1}^N R_{A_j}^{-1} \cdot t / \otimes \text{vec}[\delta_j^{-1} - U_j \cdot t //^T];$$

$$H_{m_i; m_j^T} \triangleq \frac{\partial^2 l}{\partial m_i \partial m_j^T} = -\delta_i^{-1} \sum_{t=1}^N \pm_{ij} h_i \cdot t / + \sum_{t=1}^N \circ_{ij} \cdot t / \delta_i^{-1} \cdot x \cdot t' - m_i / \times \cdot x \cdot t' - m_j /^T \delta_j^{-1};$$

$$H_{m_i; \delta_j^T} \triangleq \frac{\partial^2 l}{\partial m_i \partial \text{vec}[\delta_j]^T} = -\frac{1}{2} \sum_{t=1}^N \circ_{ij} \cdot t / \text{vec}[\delta_j^{-1} - U_j \cdot t //^T \otimes \delta_i^{-1} \cdot x \cdot t' - m_i // - \sum_{t=1}^N \pm_{ij} h_i \cdot t / \cdot \delta_i^{-1} \cdot x \cdot t' - m_i // \otimes \text{vec}[\delta_i^{-1}]^T];$$

$$H_{\delta_i; \delta_j^T} \triangleq \frac{\partial^2 l}{\partial \text{vec}[\delta_i] \partial \text{vec}[\delta_j]^T} = \frac{\partial}{\partial \text{vec}[\delta_i]} \otimes \frac{\partial}{\partial \text{vec}[\delta_j]^T} \\ = -\frac{1}{4} \sum_{t=1}^N \circ_{ij} \cdot t / \text{vec}[\delta_j^{-1} - U_j \cdot t //^T \otimes \text{vec}[\delta_i^{-1} - U_i \cdot t // \\ -\frac{1}{2} \sum_{t=1}^N \pm_{ij} h_i \cdot t / \cdot -\delta_i^{-1} \otimes \delta_i^{-1} + \delta_i^{-1} \otimes U_i \cdot t // + U_i \cdot t // \otimes \delta_i^{-1};$$

where  $\mathcal{Z} = \{\mathcal{A}; m_j; \delta_j; j = 1; \dots; K\}$ ,  $\pm_{ij}$  is the Kronecker function, and

$$H_A^{-1} \cdot t / = [h_1 \cdot t / =_{\text{O}_1}; \dots; h_K \cdot t / =_{\text{O}_K}]^T;$$

$$\circ_{ij} \cdot t / = \cdot \pm_{ij} - h_i \cdot t // h_j \cdot t /;$$

$$R_{A_j}^{-1} \cdot t / = [{}^{\circ} 1_j \cdot t / =_{\text{O}_1}; \dots; {}^{\circ} K_j \cdot t / =_{\text{O}_K}]^T;$$

$$U_i \cdot t // = U_i \cdot x \cdot t' / = \delta_i^{-1} \cdot x \cdot t' - m_i / \cdot x \cdot t' - m_i /^T \delta_i^{-1};$$

When the  $K$  gaussian distributions in the mixture are separated enough, that is, the posterior probabilities  $h_{j,t}$  are approximately zero or one at  $\mathcal{Z}^*$ , by the above expression of the Hessian matrix, it is easy to verify that  $\lim_{N \rightarrow \infty} P(\mathcal{Z}^*)/H(\mathcal{Z}^*)$  will become to the following block diagonal matrix:

$$F = \begin{pmatrix} -I_K + \mathcal{A}^*[1; 1; \dots; 1] & 0 \\ 0 & -I_{K \times (d+d^2/2)} \end{pmatrix}; \tag{2.22}$$

where  $\mathcal{A}^* = [\mathbb{1}_1^*; \dots; \mathbb{1}_K^*]^T$  and  $I_n$  is the  $n$ th order identity matrix.

Furthermore, we rearrange the  $E$  formed from the unit basis vectors that span  $\mathcal{R}_1^1$  into the following two parts:

$$E = \begin{pmatrix} E_1 & 0 \\ 0 & E_2 \end{pmatrix}; \tag{2.23}$$

with its first part  $E_1 = [e_1^1; \dots; e_{K-1}^1]$  and  $e_1^1; \dots; e_{K-1}^1$  are unit basis vectors of  $\mathcal{R}_1^1 = \{x \in \mathbb{R}^K: \sum_{j=1}^K x_j = 0\}$ .

From equations 2.22 and 2.23, we have

$$E^T \cdot I + F/E = \begin{pmatrix} E_1^T \mathcal{A}^*[1; 1; \dots; 1] E_1 & 0 \\ 0 & 0 \end{pmatrix} = 0; \tag{2.24}$$

since  $[1; \dots; 1] E_1 = 0$ .

Substituting this result into the upper bound in equation 2.21, we have:

$$r \leq \|I + E^T \lim_{N \rightarrow \infty} P(\mathcal{Z}^*)/H(\mathcal{Z}^*)/E\| = \|E^T \cdot I + F/E\| = 0;$$

Thus, the EM algorithm for gaussian mixtures has a superlinear convergence rate. A special case of this result was first given in Xu and Jordan (1996).

Actually gaussian distributions in the mixture cannot be strictly separated. In the rest of this article, we will show that  $\lim_{N \rightarrow \infty} P(\mathcal{Z}^*)/H(\mathcal{Z}^*)$  approximately converges to  $F$  as the overlap of any two gaussian distributions in the mixture tends to zero.

### 3 The Main Theorem

---

We begin by introducing a set of quantities on the overlap of component gaussian distributions as follows:

$$e_{ij}(\mathcal{Z}^*) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=1}^N |^\circ_{ij,t}| = \int |^\circ_{ij,x}| P(x|\mathcal{Z}^*) dx; \text{ for } i, j = 1; \dots; K;$$



where  $h_{ij}(x) = \frac{h_i(x)}{h_j(x)}$  and

$$h_{iL}(x) = \frac{P_i(x|m_i^*; \sigma_i^*)}{\sum_{k=1}^K P_k(x|m_k^*; \sigma_k^*)} \text{ for } L = 1; \dots; K:$$

Since  $|h_{ij}(x)| \leq 1$ ,  $e_{ij} \leq 1$ .

If  $i \neq j$ ,  $e_{ij} = \int h_{ij}(x) P(x) dx$ , which can be considered as a measure of the average overlap between the distributions of gaussian  $i$  and  $j$ . In fact, if  $P_i(x|m_i^*; \sigma_i^*)$  and  $P_j(x|m_j^*; \sigma_j^*)$  overlap in a higher degree at a point  $x$ , then  $h_{ij}(x)$  certainly takes a higher value; otherwise, it takes a lower value. When they are well separated at  $x$ ,  $h_{ij}(x)$  becomes zero. Thus, the product  $h_{ij}(x) P(x)$  represents a degree of overlap between  $P_i(x|m_i^*; \sigma_i^*)$  and  $P_j(x|m_j^*; \sigma_j^*)$  at  $x$  in the mixture environment, and the above  $e_{ij}$  is an average overlap between the distributions of gaussian  $i$  and  $j$  in the mixture.

Since two gaussian distributions always have some overlap area, we always have  $e_{ij} > 0$ . Moreover, by the fact  $\sum_{j=1}^N h_{ij}(x) = 1$ , we also have

$$\begin{aligned} e_{ii} &= \int |h_{ii}(x)| P(x) dx = \int h_{ii}(x) \cdot 1 - \sum_{j \neq i} h_{ij}(x) P(x) dx \\ &= \sum_{j \neq i} \int h_{ij}(x) P(x) dx = \sum_{j \neq i} e_{ij} > 0: \end{aligned}$$

Clearly,  $e_{ii}$  represents the overlap of the  $i$ th gaussian distribution with all the other gaussian distributions in the mixture.

We consider the worst case and define

$$e = \min_{ij} e_{ij} \leq 1: \quad (3.1)$$

By experiments, we can find that as the gaussian distributions in the mixture become more separated, the EM algorithm can asymptotically converge to the true parameter  $\theta^*$  with a better convergence rate. Here, we mathematically study the asymptotic convergence properties of

matrix  $U_i$  such that

$$\mathcal{G}_i^* = U_i^T \begin{pmatrix} \lambda_{i1} & 0 & \cdots & 0 \\ 0 & \lambda_{i2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_{id} \end{pmatrix} U_i; \tag{3.3}$$

where  $\lambda_{i1}; \dots; \lambda_{id}$  are the eigenvalues of  $\mathcal{G}_i^*$ . By taking the transformation  $y = U_i \cdot x - m_i^*$ , we have from equation 3.2 that

$$\sum_{j=1}^n \lambda_{ij}^{-1} y_j^2 = 1;$$

which is a hyperellipsoid at the origin in the transformed space. It follows from the inverse transformation  $x = m_i^* + U_i^T y$  that the characteristic contour is a hyperellipsoid at  $m_i^*$  in the original space. In fact, the axes of this characteristic hyperellipsoid describe the attenuating rates of the component density in the coordinate basis directions in the transformed space (or the  $d$  principal directions in the original space from  $m_i^*$ ), respectively, since they are the square roots of the eigenvalues, which are actually the variances of the  $i$ th component distribution in these directions, respectively. In other words, the density will attenuate rapidly as a point goes away from the center along a basis direction if the corresponding axis is small; otherwise, it will attenuate slowly. For simplicity, we use the minimum-radius hypersphere, which includes the characteristic contour,

$$\|x - m_i^*\|^T \cdot \|x - m_i^*\| = \lambda_{i \max}^i;$$

instead of the characteristic contour, and then the radius of this hypersphere is  $\lambda_{i \max}^i / \sqrt{2}$ . In the same way, we can have the characteristic contour and the mini-radius hypersphere for the  $j$ th component distribution. Since the overlap between these two densities is approximately proportional to the overlap between the two characteristic contours or minimum-radius hyperspheres, we can define  $\lambda_{ij}^*$  as the overlap between the densities of components  $i$  and  $j$  as follows:

$$\lambda_{ij}^* = \frac{\lambda_{i \max}^i / \sqrt{2} \cdot \lambda_{j \max}^j / \sqrt{2}}{\|m_i^* - m_j^*\|};$$

Therefore, for a gaussian mixture of true parameter  $2^*$ , we define

$$\lambda^* = \max_{i \neq j} \lambda_{ij}^* \tag{3.4}$$

as an overall overlap degree of gaussians in the mixture from point of view of the worst case. Obviously,  $\epsilon \cdot 2^* \rightarrow 0$  is equivalent to  $e \cdot 2^* \rightarrow 0$ .

**Second**, we consider some assumptions that regularize the manner of  $e \cdot 2^*$  tending to zero.

We first assume that  $2^*$  satisfies

$$\text{Condition 1: } \alpha_i^* \geq \alpha; \text{ for } i = 1; \dots; K;$$

where  $\alpha$  is a positive number. When some mixing proportion  $\alpha_i^*$  tends to zero, the corresponding gaussian will withdraw from the mixture, which degenerates to a mixture with a lower number of the mixing components. This assumption excludes this degeneracy.

Our second assumption is that the eigenvalues of all the covariance matrices satisfy

$$\text{Condition 2: } \tau \cdot 2^* \leq \lambda_{i,k} \leq \nu \cdot 2^*; \text{ for } i = 1; \dots; K; \quad k = 1; \dots; d;$$

where  $\tau$  is a positive number and  $\nu \cdot 2^*$  is defined to be the maximum eigenvalue of the covariance matrices  $\Sigma_1^*; \dots; \Sigma_K^*$ , that is,

$$\nu \cdot 2^* = \max_{i,k} \lambda_{i,k};$$

which is always upper bounded by a positive number  $B$ . That is, all the eigenvalues uniformly attenuate or reduce to zero when they tend to zero. It follows from condition 2 that condition numbers of all the covariance matrices are uniformly upper bounded, that is,

$$1 \leq \kappa \cdot \Sigma_i^* \leq B'; \text{ for } i = 1; \dots; K;$$

where  $\kappa \cdot \Sigma_i^*$  is the condition number of  $\Sigma_i^*$  and  $B'$  is a positive number.

The third assumption is that the mean vectors of the component densities in the mixture satisfy

$$\begin{aligned} \text{Condition 3: } \alpha D_{\max} \cdot 2^* \leq D_{\min} \cdot 2^* \leq \|m_i^* - m_j^*\| \\ \leq D_{\max} \cdot 2^*; \text{ for } i \neq j; \end{aligned}$$

where  $D_{\max} \cdot 2^* = \max_{i \neq j} \|m_i^* - m_j^*\|$ ;  $D_{\min} \cdot 2^* = \min_{i \neq j} \|m_i^* - m_j^*\|$ , and  $\alpha$  is a positive number. That is, all the distances between two mean vectors are the same-order in nitely large quantities when they tend to infinity. Moreover, when the overlap of densities in the mixture reduces to zero, any of two means  $m_i^*$ ;  $m_j^*$  cannot be arbitrarily close; there should be a positive value  $T$  such that  $\|m_i^* - m_j^*\| \geq T$  when  $i \neq j$ . Also, it is natural to assume that the mean vectors take different directions when they diverge to infinity.

With the above preparations, we are ready to introduce our main theorem.

**Theorem 1.** Given i.i.d. samples  $\{x^t\}_1^N$  from a mixture of  $K$  gaussian distributions of the parameters  $\mathcal{Z}^*$  that satisfies conditions 1–3, when  $e \cdot \mathcal{Z}^*/$  is considered as an infinitesimal, as it tends to zero, we have:

$$\lim_{N \rightarrow \infty} P \cdot \mathcal{Z}^*/H \cdot \mathcal{Z}^*/ = F + o \cdot e^{0:5-\epsilon} \cdot \mathcal{Z}^*//; \tag{3.5}$$

where  $F$  is given by equation 2.22 and  $\epsilon$  is an arbitrarily small positive number.

According to this theorem, as the overlap of distributions in the mixture becomes small or, more precisely,  $e \cdot \mathcal{Z}^*/ \rightarrow 0$ , the asymptotic value of  $P \cdot \mathcal{Z}^*/H \cdot \mathcal{Z}^*/$  becomes  $F$  plus a matrix in which each element is a higher-order infinitesimal of  $e^{0:5-\epsilon} \cdot \mathcal{Z}^*/$ . It further follows from equation 2.24 that  $E^T \cdot I + F/E = 0$ . Thus, we have

$$\begin{aligned} I + E^T \lim_{N \rightarrow \infty} P \cdot \mathcal{Z}^*/H \cdot \mathcal{Z}^*/E &= E^T \cdot I + \lim_{N \rightarrow \infty} P \cdot \mathcal{Z}^*/H \cdot \mathcal{Z}^*//E \\ &= o \cdot e^{0:5-\epsilon} \cdot \mathcal{Z}^*//; \end{aligned}$$

That is, each element of  $I + E^T \lim_{N \rightarrow \infty} P \cdot \mathcal{Z}^*/H \cdot \mathcal{Z}^*/E$  is a higher-order infinitesimal of  $e^{0:5-\epsilon} \cdot \mathcal{Z}^*/$ . Because the norm of a real matrix  $A = [a_{ij}]_{m \times m}$  is always not larger than  $\sqrt{m} \max_i \sum_{j=1}^m |a_{ij}|$ , the norm  $\|I + E^T P \cdot \mathcal{Z}^*/H \cdot \mathcal{Z}^*/E\|$  is certainly a higher-order infinitesimal of  $e^{0:5-\epsilon} \cdot \mathcal{Z}^*/$ . Therefore, it follows from equation 2.21 that the asymptotic convergence rate of the EM algorithm locally around  $\mathcal{Z}^*$  is a higher-order infinitesimal of  $e^{0:5-\epsilon} \cdot \mathcal{Z}^*/$ . That is, when  $e \cdot \mathcal{Z}^*/$  is small and  $N$  is large enough, the convergence rate of the EM algorithm approaches approximately zero. In other words, the EM algorithm in this case has a quasi-Newton-type convergence behavior, with its asymptotic convergence rate dominated by the infinitesimal  $e^{0:5-\epsilon} \cdot \mathcal{Z}^*/$ .

From this theorem, we can also find that  $P \cdot \mathcal{Z}^*/$  tends to the inverse of  $H \cdot \mathcal{Z}^*/$  as  $e \cdot \mathcal{Z}^*/$  tends to zero. Thus, when the overlap of gaussian distributions in the mixture becomes a small value, the EM algorithm approximates to the Newton algorithm. Actually,  $P \cdot \mathcal{Z}^*/$  makes the EM algorithm for the gaussian mixture be different from the conventional first-order iterative algorithm such as gradient ascent. This may be the basic underlying factor that speeds up the convergence of the EM algorithm when the overlap of gaussians in the mixture is low.

#### 4 Lemmas

---

In this section, we describe several lemmas that will be used for proving the main theorem. The norm for a vector or matrix is always assumed to be Euclidean norm in the following.

We now define three kinds of special polynomial functions that we often meet in the further analyses:

**Definition 1.**  $g(x; \theta^*)$  is called a regular function if it satisfies:

- (i) If  $\theta^*$  is fixed,  $g(x; \theta^*)$  is a polynomial function of the component variables  $x_1; \dots; x_d$  of  $x$ .
- (ii) If  $x$  is fixed,  $g(x; \theta^*)$  is a polynomial function of the elements of  $m_1^*; \dots; m_K^*, \phi_1^*; \dots; \phi^*$

As  $E. \|Y\|^i |0; I_d/$  is certainly finite and  $\dots 2^*/$  is upper bounded,  $E. \|X - m_j^* \|^i |m_j^*; \delta_j^*/$  is upper bounded. Since  $\theta_j^* < 1$  and  $m. 2^*/$  is lower bounded by  $\frac{T}{2}$  under the constraint  $\|m_i^* - m_j^*\| \geq T$ , it follows from equation 4.2 that there exist a positive number  $\epsilon^1$  and a positive integer  $q$  such that

$$E. \|X\|^k |2^*/ = \int \|x\|^k P.x |2^*/ dx \leq \epsilon^1 m^q . 2^*/: \tag{4.4}$$

We further get an upper bound for  $E. g^2 . X; 2^*/$ . In the process of  $e. 2^*/$  tending to zero under the assumptions, the elements of  $\mathcal{A}^*$  and  $\mathcal{A}^{*-1}$  are bounded. Since  $\|\delta_j^*\| = \dots_j \max \leq \dots . 2^*/ \leq B$ , the elements of each  $\delta_j^*$  are bounded. Moreover, since  $\dots_{j1}; \dots; \dots_{jd}$  are the eigenvalues of  $\delta_j^*$  and  $\dots_{j1}^{-1}; \dots; \dots_{jd}^{-1}$  are the eigenvalues of  $\delta_j^{*-1}$ , we have

$$\|\dots . 2^*/ \delta_j^{*-1}\| \leq \dots . 2^*/ \frac{1}{\dots . 2^*/} = \frac{1}{\dots}:$$

That is, the elements of  $\dots . 2^*/ \delta_j^{*-1}$  are also bounded. We consider that  $g.x; 2^*/$  is a polynomial function of  $x_1; \dots; x_d$  and that the coefficient of each term in  $g.x; 2^*/$  is a polynomial function of the elements of  $\mathcal{A}^*, \mathcal{A}^{*-1}, m_1^*; \dots; m_K^*, \delta_1^*; \dots; \delta_K^*, \dots . 2^*/ \delta_1^{*-1}; \dots; \dots . 2^*/ \delta_K^{*-1}$  with constant coefficients. Because the elements of these matrices and vectors except  $m_1^*; \dots; m_K^*$  are bounded, the absolute value of the coefficient is certainly upper bounded by a positive-order power function of  $m. 2^*/$ . Therefore,  $|g.x; 2^*/|$  is upper bounded by a positive polynomial function of  $\|x\|$  with its coefficients being some polynomial functions of  $m. 2^*/$  of positive constant coefficients. As  $g^2.x; 2^*/$  is certainly a balanced function, it is also upper bounded by a positive polynomial function of  $\|x\|$ . According to equation 4.4 and the property that  $m. 2^*/$  is lower bounded by  $\frac{T}{2}$ , there certainly exist a positive number  $\epsilon^1$  and a positive integer  $p$  such that

$$E. g^2 . X; 2^*/ \leq \epsilon^1 m^{2p} . 2^*/:$$

Finally, it follows from the Cauchy-Schwarz inequality and the fact  $|\epsilon_{ij}.x| \leq 1$  that

$$\begin{aligned} \int |g.x; 2^*/ \epsilon_{ij}.x| |P.x |2^*/ dx &= E. |g.X; 2^*/| |\epsilon_{ij}.X| / \\ &\leq E^{\frac{1}{2}} . g^2 . X; 2^*/ / E^{\frac{1}{2}} . |\epsilon_{ij}.X| / \\ &\leq \epsilon^1 m^p . 2^*/ / \sqrt{\epsilon_{ij}. 2^*/} \\ &\leq \epsilon^1 m^p . 2^*/ / \sqrt{\epsilon . 2^*/}: \end{aligned}$$

Since  $\hat{e}(\theta^*)$  is not an invertible function, that is, there may be many  $\theta^*$  for a value of  $\hat{e}(\theta^*)$ , we further define

$$f(\hat{e}) = \sup_{\theta^* / \hat{e}(\theta^*)} e(\theta^*) \tag{4.5}$$

Because  $e(\theta^*)$  is always not larger than 1 by the definition,  $f(\hat{e})$  is well defined. By this definition, we certainly have

$$e_{ij}(\theta^*) \leq e(\theta^*) \leq f(\hat{e}(\theta^*)) \tag{4.6}$$

**Lemma 2.** Suppose that  $\theta^*$  satisfies conditions 1–3. As  $\hat{e}(\theta^*)$  tends to zero, we have:

(i)  $\hat{e}(\theta^*)$ ,  $\hat{e}_{ij}(\theta^*)$  and  $\frac{1}{\|m_i^* - m_j^*\|}$  are the equivalent infinitesimals.

(ii) For  $i \neq j$ , we have

$$\|m_i^*\| \leq T' \|m_i^* - m_j^*\| \tag{4.7}$$

where  $T'$  is a positive number.

(iii) For any two nonnegative numbers  $p$  and  $q$  with  $p + q > 0$ , we have

$$\|m_i^* - m_j^*\|^p \cdot \hat{e}_{\max}^{-q} \leq O(\hat{e}^{-p \vee q}) \tag{4.8}$$

$$\|m_i^* - m_j^*\|^p \cdot \hat{e}_{\max}^{-q} \geq O(\hat{e}^{-p \wedge q}) \tag{4.9}$$

where  $p \vee q = \max\{p, q\}$ ;  $p \wedge q = \min\{p, q\}$ .

See the appendix for the proof.

**Lemma 3.** When  $\theta^*$  satisfies conditions 1–3 and  $\hat{e}(\theta^*) \rightarrow 0$  as an infinitesimal, we have

$$f(\hat{e}(\theta^*)) = o(\hat{e}^p) \tag{4.10}$$

where  $p > 0$ ,  $p$  is any positive number and  $o(x)$  means that it is a higher-order infinitesimal as  $x \rightarrow 0$ .

See the appendix for the proof.

**Lemma 4.** Suppose that  $g(x; \theta^*)$  is a regular and convertible function and that  $\theta^*$  satisfies conditions 1–3. As  $e(\theta^*) \rightarrow 0$  is considered as an infinitesimal, we have

$$\begin{aligned} \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=1}^N \int_{\theta^*} g(x; \theta^*) x^t; \theta^* &= \int g(x; \theta^*) f_{ij}(x) P(x | \theta^*) dx \\ &= o(e^{0.5 - \epsilon}) \end{aligned} \tag{4.11}$$

where  $\epsilon$  is an arbitrarily small, positive number.





## 5 Proof of Theorem 1

---

We are now ready to prove the main theorem.

**Proof of Theorem 1.** We are interested in the matrix  $P \cdot 2^*/H \cdot 2^*$ , which determines the local convergence rate of EM.

The explicit expressions of  $P \cdot 2^*$  and  $H \cdot 2^*$  allow us to obtain formulas for the blocks of  $P \cdot 2^*/H \cdot 2^*$ . To clarify our notation, we begin by writing out these blocks as follows:

$$P \cdot 2^*/H \cdot 2^*$$

$$= \text{diag}[P_{\mathcal{A}}; P_{m_1}; \dots; P_{m_K}; P_{\delta_1}; \dots; P_{\delta_K}]$$

$$\times \begin{pmatrix} H_{\mathcal{A};\mathcal{A}^T} & H_{\mathcal{A};m_1^T} & \cdots & H_{\mathcal{A};m_K^T} & H_{\mathcal{A};\delta_1^T} & \cdots & H_{\mathcal{A};\delta_K^T} \\ H_{m_1;\mathcal{A}^T} & H_{m_1;m_1^T} & \cdots & H_{m_1;m_K^T} & H_{m_1;\delta_1^T} & \cdots & H_{m_1;\delta_K^T} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ H_{m_K;\mathcal{A}^T} & H_{m_K;m_1^T} & \cdots & H_{m_K;m_K^T} & H_{m_K;\delta_1^T} & \cdots & H_{m_K;\delta_K^T} \\ H_{\delta_1;\mathcal{A}^T} & H_{\delta_1;m_1^T} & \cdots & H_{\delta_1;m_K^T} & H_{\delta_1;\delta_1^T} & \cdots & H_{\delta_1;\delta_K^T} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ H_{\delta_K;\mathcal{A}^T} & H_{\delta_K;m_1^T} & \cdots & H_{\delta_K;m_K^T} & H_{\delta_K;\delta_1^T} & \cdots & H_{\delta_K;\delta_K^T} \end{pmatrix}$$

$$= \begin{pmatrix} P_{\mathcal{A}}H_{\mathcal{A};\mathcal{A}^T} & P_{\mathcal{A}}H_{\mathcal{A};m_1^T} & \cdots & P_{\mathcal{A}}H_{\mathcal{A};m_K^T} & P_{\mathcal{A}}H_{\mathcal{A};\delta_1^T} & \cdots & P_{\mathcal{A}}H_{\mathcal{A};\delta_K^T} \\ P_{m_1}H_{m_1;\mathcal{A}^T} & P_{m_1}H_{m_1;m_1^T} & \cdots & P_{m_1}H_{m_1;m_K^T} & P_{m_1}H_{m_1;\delta_1^T} & \cdots & P_{m_1}H_{m_1;\delta_K^T} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ P_{m_K}H_{m_K;\mathcal{A}^T} & P_{m_K}H_{m_K;m_1^T} & \cdots & P_{m_K}H_{m_K;m_K^T} & P_{m_K}H_{m_K;\delta_1^T} & \cdots & P_{m_K}H_{m_K;\delta_K^T} \\ P_{\delta_1}H_{\delta_1;\mathcal{A}^T} & P_{\delta_1}H_{\delta_1;m_1^T} & \cdots & P_{\delta_1}H_{\delta_1;m_K^T} & P_{\delta_1}H_{\delta_1;\delta_1^T} & \cdots & P_{\delta_1}H_{\delta_1;\delta_K^T} \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ P_{\delta_K}H_{\delta_K;\mathcal{A}^T} & P_{\delta_K}H_{\delta_K;m_1^T} & \cdots & P_{\delta_K}H_{\delta_K;m_K^T} & P_{\delta_K}H_{\delta_K;\delta_1^T} & \cdots & P_{\delta_K}H_{\delta_K;\delta_K^T} \end{pmatrix} :$$

Based on the expressions of the Hessian blocks, we have:

$$\begin{aligned} P_{\mathcal{A}}H_{\mathcal{A};\mathcal{A}^T} &= \frac{1}{N} \cdot \text{diag}[\otimes_1^*; \dots; \otimes_K^*] - \mathcal{A}^* \cdot \mathcal{A}^{*/T} / \left( - \sum_{t=1}^N H_A^{-1} \cdot t \cdot H_A^{-1} \cdot t / T \right) \\ &= - \cdot \text{diag}[\otimes_1^*; \dots; \otimes_K^*] - \mathcal{A}^* \cdot \mathcal{A}^{*/T} / \left( \frac{1}{N} \sum_{t=1}^N H_A^{-1} \cdot t \cdot H_A^{-1} \cdot t / T \right) : \end{aligned}$$

By taking the limit as  $N \rightarrow \infty$ , we obtain:

$$\lim_{N \rightarrow \infty} P_{\mathcal{A}} H_{\mathcal{A}} : \mathcal{A}^T = - \text{diag}[\otimes_1^* ; \dots ; \otimes_K^*] - \mathcal{A}^* \cdot \mathcal{A}^{*T} / \lim_{N \rightarrow \infty} \frac{1}{N} \\ \times \sum_{t=1}^N H_A^{-1} \cdot t / \cdot H_A^{-1} \cdot t //^T :$$

We further let

$$I \cdot 2^* / = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=1}^N H_A^{-1} \cdot t / H_A^{-1} \cdot t //^T \tag{5.1}$$

and compute this matrix by the elements.

When  $i \neq j$ , we have

$$I \cdot 2^* / \cdot i ; j / = \frac{1}{\otimes_1^* \otimes_j^*} \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=1}^N h_i \cdot t / h_j \cdot t / = \frac{1}{\otimes_i^* \otimes_j^*} e_{ij} \cdot 2^* / \\ \leq \frac{1}{\otimes_2} e \cdot 2^* / = 0 \cdot e^{0.5 - \dots} \cdot 2^* // :$$

When  $i = j$ , we further have

$$I \cdot 2^* / \cdot i ; i / = \frac{1}{\otimes_i^* / 2} \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=1}^N h_i \cdot t / [1 - \cdot 1 - h_i \cdot t /] \\ = \frac{1}{\otimes_i^* / 2} \int h_i \cdot x / P \cdot x | 2^* / dx - \frac{1}{\otimes_i^* / 2} \int | \circ_{ii} \cdot x / | P \cdot x | 2^* / dx \\ \geq \otimes_i^{*-1} - \frac{1}{\otimes_2} e \cdot 2^* / = \otimes_i^{*-1} + 0 \cdot e^{0.5 - \dots} \cdot 2^* // :$$

Then we have

$$I \cdot 2^* / = \text{diag}[\otimes_1^{*-1} ; \dots ; \otimes_K^{*-1}] + 0 \cdot e^{0.5 - \dots} \cdot 2^* // : \tag{5.2}$$

Therefore, we get

$$\lim_{N \rightarrow \infty} P_{\mathcal{A}} H_{\mathcal{A}} : \mathcal{A}^T \\ = - \text{diag}[\otimes_1^* ; \dots ; \otimes_K^*] - \mathcal{A}^* \cdot \mathcal{A}^{*T} / I \cdot 2^* / \\ = - \text{diag}[\otimes_1^* ; \dots ; \otimes_K^*] - \mathcal{A}^* \cdot \mathcal{A}^{*T} / \cdot \text{diag}[\otimes_1^{*-1} ; \dots ; \otimes_K^{*-1}] + 0 \cdot e^{0.5 - \dots} \cdot 2^* // \\ = -I_K + \mathcal{A}^* [1 ; 1 ; \dots ; 1] + 0 \cdot e^{0.5 - \dots} \cdot 2^* // :$$

For  $j = 1; \dots; K$ , we have

$$\begin{aligned} P_{\mathcal{A}} H_{\mathcal{A}; m_j^T} &= \frac{1}{N} \cdot \text{diag}[\theta_1^*; \dots; \theta_K^*] - \mathcal{A}^* \cdot \mathcal{A}^{*T} / \sum_{t=1}^N R_{A_j}^{-1} \cdot t / x^{t/} - m_j^{*T} / \theta_j^{*-1} \\ &= \text{diag}[\theta_1^*; \dots; \theta_K^*] - \mathcal{A}^* \cdot \mathcal{A}^{*T} / \frac{1}{N} \sum_{t=1}^N R_{A_j}^{-1} \cdot t / x^{t/} - m_j^{*T} / \theta_j^{*-1} \end{aligned}$$

Since  $R_{A_j}^{-1} \cdot t / = [\theta_1^*; \dots; \theta_K^*]^T$ , we let  $g; x; \theta_j^{*-1}$  be any element of  $\mathcal{A}^{*-1} \cdot x - m_j^{*T} / \theta_j^{*-1}$ . Obviously, this kind of  $g; x; \theta_j^{*-1}$  is a regular and convertible function with  $q = 1$ . By lemma 4, we have:

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=1}^N R_{A_j}^{-1} \cdot t / x^{t/} - m_j^{*T} / \theta_j^{*-1} = 0 \cdot e^{0.5 - \dots} \cdot \theta_j^{*-1} \quad (5.3)$$

Because the elements of the matrix  $\text{diag}[\theta_1^*; \dots; \theta_K^*] - \mathcal{A}^* \cdot \mathcal{A}^{*T} /$  are bounded, we further have

$$\lim_{N \rightarrow \infty} P_{\mathcal{A}} H_{\mathcal{A}; m_j^T} = 0 \cdot e^{0.5 - \dots} \cdot \theta_j^{*-1} \quad (5.4)$$

Because the Hessian matrix has the following property,

$$H_{m_j; \mathcal{A}}$$

Next, we consider  $P_{m_i} H_{m_i; m_j^T}$  and have

$$P_{m_i} H_{m_i; m_j^T} = -\pm_{ij} I_d + \quad 1$$

On the other hand, we have

$$\begin{aligned}
 P_{\delta_i} H_{\delta_i; m_j^T} &= \frac{2}{\sum_{t=1}^N h_{i,t}} \delta_i^* \otimes \delta_i^* H_{m_j; \delta_i^T}^T \\
 &= -\delta_i^* \otimes \delta_i^* \left\{ \frac{1}{\sum_{t=1}^N h_{i,t}} \sum_{t=1}^N \circ_{ij,t} / \text{vec}[\delta_i^{*-1} - U_{i,t}] \right. \\
 &\quad \left. \otimes \{\delta_j^{*-1} \cdot x^{t/} - m_j^* / \}^T \right\} \\
 &\quad - 2\delta_i^* \otimes \delta_i^* \left\{ \frac{1}{\sum_{t=1}^N h_{i,t}} \sum_{t=1}^N \pm_{ij} h_{i,t} / \{\delta_i^{*-1} \cdot x^{t/} - m_i^* / \}^T \right. \\
 &\quad \left. \otimes \text{vec}[\delta_i^{*-1}] \right\};
 \end{aligned}$$

and by the same argument used for  $P_{m_i} H_{m_i; \delta_j^T}$ , we have

$$\lim_{N \rightarrow \infty} P_{\delta_i} H_{\delta_i; m_i^T} = 0.e^{0.5-} \cdot 2^{*//}; \quad (5.10)$$

Now we turn to the remaining case,  $P_{\delta_i} H_{\delta_i; \delta_j^T}$ , and have:

$$\begin{aligned}
 P_{\delta_i} H_{\delta_i; \delta_j^T} &= -\frac{1}{2} \delta_i^* \otimes \delta_i^* \left\{ \frac{1}{\sum_{t=1}^N h_{i,t}} \sum_{t=1}^N \circ_{ij,t} / \text{vec}[\delta_j^{*-1} - U_{j,t}] \right. \\
 &\quad \left. \otimes \text{vec}[\delta_i^{*-1} - U_{i,t}] \right\} \\
 &\quad - \pm_{ij} \cdot \delta_i^* \otimes \delta_i^* / \{-\delta_i^{*-1} \otimes \delta_i^{*-1} + \delta_i^{*-1} \\
 &\quad \otimes \left\{ \frac{1}{\sum_{t=1}^N h_{i,t}} \sum_{t=1}^N h_{i,t} / U_{i,t} \right\} \\
 &\quad + \left\{ \frac{1}{\sum_{t=1}^N h_{i,t}} \sum_{t=1}^N h_{i,t} / U_{i,t} \right\} \otimes \delta_i^{*-1} \};
 \end{aligned}$$

By corollary 1 and the law of large number, we have

$$\begin{aligned}
 \lim_{N \rightarrow \infty} P_{\delta_i} H_{\delta_i; \delta_j^T} &= 0.e^{0.5-} \cdot 2^{*//} - \pm_{ij} \cdot \delta_i^* \otimes \delta_i^* / \cdot \delta_i^{*-1} \otimes \delta_i^{*-1} / \\
 &= -\pm_{ij} I_{d^2} + 0.e^{0.5-} \cdot 2^{*//};
 \end{aligned}$$

Summing up all the results, we obtain:

$$\lim_{N \rightarrow \infty} \frac{P(\cdot | \mathcal{H}_N)}{P(\cdot | \mathcal{H}_N)} = \begin{pmatrix} 0 & -I_{K \times (d+d^2)} \end{pmatrix} + o(e^{-0.5N})$$

$$= F + o(e^{-0.5N});$$

and the proof is completed.

### 6 Conclusions

We have presented an analysis of the asymptotic convergence rate of the EM algorithm for gaussian mixtures. Our analysis shows that when the overlap of any two gaussian distributions is small enough, the large sample local convergence behavior of the EM algorithm is regular. (large)

$\|m_i^*\|$  or  $\|m_j^*\|$  is upper bounded. Otherwise, if  $\|m_i^*\|$  and  $\|m_j^*\|$  both increase to infinity, since  $m_i^*$  and  $m_j^*$  take different directions, the order of the infinitely large quantity  $\|m_i^*\|$  is certainly lower than or equal to that of  $\|m_i^* - m_j^*\|$ , which also leads to equation 4.7. Therefore, (ii) holds under the assumptions.

Finally, we turn to (iii). For the first inequality, equation 4.8, we consider three cases, as follows.

In the simple case  $p = q > 0$ , according to (i), we have

$$\begin{aligned} \|m_i^* - m_j^*\|^p \cdot \rho_{\max}^{i/q} &= \|m_i^* - m_j^*\|^p \cdot \rho_{\max}^{i/q} \\ &= \left( \frac{\rho_{\max}^i}{\|m_i^* - m_j^*\|} \right)^{-p} = O(\rho_{\max}^{-p} \cdot 2^{*//}) = O(\rho_{\max}^{-p \vee q} \cdot 2^{*//}); \end{aligned}$$

If  $p > q$ , since  $\rho_{\max}^i$  is upper bounded and according to (i), we have

$$\|m_i^* - m_j^*\|^p \cdot \rho_{\max}^{i/q} \leq O(\rho_{\max}^{-p} \cdot 2^{*//}) = O(\rho_{\max}^{-p \vee q} \cdot 2^{*//});$$

If  $p < q$ , as  $\|m_i^* - m_j^*\| \geq T$ , we can have

$$\|m_i^* - m_j^*\|^p \cdot \rho_{\max}^{i/q} \leq O(\rho_{\max}^{-q} \cdot 2^{*//}) = O(\rho_{\max}^{-p \vee q} \cdot 2^{*//});$$

Summing up the results on the three cases, we obtain:

$$\|m_i^* - m_j^*\|^p \cdot \rho_{\max}^{i/q} \leq O(\rho_{\max}^{-p \vee q} \cdot 2^{*//});$$

In a similar way as above, we can prove the second inequality, equation 4.9.

**Proof of Lemma 3.** We first prove that

$$f(\rho) = o(\rho^p);$$

as  $\rho \rightarrow 0$ , where  $p$  is an arbitrarily positive number.

We consider the mixture of  $K$  gaussian densities of the parameter  $2^*$  under the relation  $\rho \cdot 2^* / = \rho$ . When  $i \neq j$ , for a small enough  $\rho$ , there is certainly a point  $m_{ij}^*$  on the line between  $m_i^*$  and  $m_j^*$  such that

$$\rho_i^* P_i \cdot m_{ij}^* | m_i^*; \rho_i^* / = \rho_j^* P_j \cdot m_{ij}^* | m_j^*; \rho_j^* /;$$

We further define

$$E_i = \{x: \rho_i^* P_i \cdot x | m_i^*; \rho_i^* / \geq \rho_j^* P_j \cdot x | m_j^*; \rho_j^* / \};$$

$$E_j = \{x: \rho_j^* P_j \cdot x | m_j^*; \rho_j^* / > \rho_i^* P_i \cdot x | m_i^*; \rho_i^* / \};$$

As  $\epsilon$  tends to zero,  $\frac{\epsilon}{\|m_i^* - m_j^*\|}$  and  $\frac{\epsilon}{\|m_i^* - m_j^*\|}$  are the same-order infinitesimals. Moreover,  $\epsilon$  and  $\epsilon$  are both upper bounded. Thus, there certainly exist some neighborhoods of  $m_i^*$  (or  $m_j^*$ ) in  $E_i$  (or  $E_j$ ). For clarity, we let  $\mathcal{N}_{r_i, m_i^*}$  and  $\mathcal{N}_{r_j, m_j^*}$  be the largest neighborhood in  $E_i$  and  $E_j$ , respectively, where  $r_i$  and  $r_j$  are their radii. Obviously  $r_i$  and  $r_j$  are both proportional to  $\|m_i^* - m_j^*\|$  when  $\|m_i^* - m_j^*\|$  either tends to infinity or is upper bounded. So there exist a pair of positive numbers  $b_1$  and  $b_2$  such that

$$r_i \geq b_1 \|m_i^* - m_j^*\|; \text{ and } r_j \geq b_2 \|m_i^* - m_j^*\|;$$

We further define

$$\begin{aligned} \mathcal{D}_i &= \mathcal{N}_{r_i, m_i^*}^c = \{x: \|x - m_i^*\| \geq r_i\}; \\ \mathcal{D}_j &= \mathcal{N}_{r_j, m_j^*}^c = \{x: \|x - m_j^*\| \geq r_j\}; \end{aligned}$$

and thus

$$E_i \subset \mathcal{D}_j; \quad E_j \subset \mathcal{D}_i;$$

Moreover, from the definitions of  $e_{ij}$  and  $h_k(x)$ , we have

$$\begin{aligned} e_{ij} &= \int h_i(x)/h_j(x) P(x) dx \\ &= \int_{E_i} h_i(x)/h_j(x) P(x) dx + \int_{E_j} h_i(x)/h_j(x) P(x) dx \\ &\leq \int_{\mathcal{D}_j} h_i(x)/h_j(x) P(x) dx + \int_{\mathcal{D}_i} h_i(x)/h_j(x) P(x) dx \\ &= \int_{\mathcal{D}_j} P_j(x|m_j^*; \mathcal{D}_j^*) dx + \int_{\mathcal{D}_i} P_i(x|m_i^*; \mathcal{D}_i^*) dx; \end{aligned}$$

We now consider  $\int_{\mathcal{D}_i} P_i(x|m_i^*; \mathcal{D}_i^*) dx$ . Since  $r_i \geq b_1 \|m_i^* - m_j^*\|$ ,

$$\int_{\mathcal{D}_i} P_i(x|m_i^*; \mathcal{D}_i^*) dx \leq \int_{\|x - m_i^*\| \geq b_1 \|m_i^* - m_j^*\|} P_i(x|m_i^*; \mathcal{D}_i^*) dx;$$

When  $\epsilon$  is sufficiently small, by the transformation  $y = U_i(x - m_i^*)/\epsilon$



$m_j^*$ , we have

$$\int_{\mathcal{D}_i} P_i(x|m_i^*; \mathcal{G}_i^*) dx \leq \int_{\|y\| \geq b_i} \frac{\|m_i^* - m_j^*\|}{i_{\max}^{d/2}} e^{-\frac{1}{2} \frac{\|m_i^* - m_j^*\|^2}{i_{\max}^d} \|y\|^2} dy; \tag{A.4}$$

where  $\circ$  is a positive number.

By lemma 2, we have

$$\|m_i^* - m_j^*\| \cdot i_{\max}^{-d/2} \leq O(\rho^{-c_1} \cdot 2^{*//});$$

$$\|m_i^* - m_j^*\|^2 \cdot i_{\max}^{-d} \geq O(\rho^{-c_2} \cdot 2^{*//});$$

where  $c_1 = 1 \vee \frac{d}{2}$ ;  $c_2 = 2 \wedge d$ .

According to these results, we have from equation A.4 that

$$\int_{\mathcal{D}_i} P_i(x|m_i^*; \mathcal{G}_i^*) dx \leq \int_{\mathcal{B}_i} \frac{1}{\rho^{c_1}} e^{-\frac{1}{2} \frac{1}{c_2} \|y\|^2} dy; \tag{A.5}$$

where  $\mathcal{B}_i = \{y: \|y\| \geq b_i\}$ ,  $1$  and  $\frac{1}{2}$  are positive numbers.

Furthermore, we let

$$F_i(\rho) = \int_{\mathcal{B}_i} P_i(y|\rho) dy; \quad P_i(y|\rho) = \frac{1}{\rho^{c_1}} e^{-\frac{1}{2} \frac{1}{c_2} \|y\|^2};$$

and consider the limit of  $\frac{F_i(\rho)}{\rho^p}$  as  $\rho$  tends to zero.

For each  $y \in \mathcal{B}_i$ , we have

$$\lim_{\rho \rightarrow 0} \frac{P_i(y|\rho)}{\rho^p} = \lim_{\rho \rightarrow 0} \frac{\rho^{c_1+p}}{e^{c_2/2 \|y\|^2}} = 0$$

uniformly in  $\mathcal{B}_i$ , which leads to

$$\lim_{\rho \rightarrow 0} \frac{F_i(\rho)}{\rho^p} = \int_{\mathcal{B}_i} \lim_{\rho \rightarrow 0} \frac{P_i(y|\rho)}{\rho^p} dy = 0;$$

and thus  $F_i(\rho) = o(\rho^p)$ : It further follows from equation A.5 that

$$\sup_{\rho \in (0, \rho^*)} \int_{\mathcal{D}_i} P_i(x|m_i^*; \mathcal{G}_i^*) dx = o(\rho^p); \tag{A.6}$$

Similarly, we can also prove

$$\sup_{\rho \in (0, \rho^*)} \int_{\mathcal{D}_j} P_j(x|m_j^*; \mathcal{G}_j^*) dx = o(\rho^p);$$

As a result, we have

$$\begin{aligned}
 f_{ij}(\cdot) &= \sup_{\mathcal{D}_j} e_{ij}(\cdot) \\
 &\leq \sup_{\mathcal{D}_j} \left( \int_{\mathcal{D}_j} P_j(x) |m_j^*(x); \delta_j^*(x)| dx + \int_{\mathcal{D}_i} P_i(x) |m_i^*(x); \delta_i^*(x)| dx \right) \\
 &\leq \sup_{\mathcal{D}_j} \int_{\mathcal{D}_j} P_j(x) |m_j^*(x); \delta_j^*(x)| dx + \sup_{\mathcal{D}_i} \int_{\mathcal{D}_i} P_i(x) |m_i^*(x); \delta_i^*(x)| dx \\
 &= 0.
 \end{aligned}$$

In the case  $i = j$ , we also have

$$\begin{aligned}
 f_{ii}(\cdot) &= \sup_{\mathcal{D}_i} \int |i_{ii}(x) / P(x)|^2 dx \\
 &= \sup_{\mathcal{D}_i} \int h_i(x) / (1 - h_i(x) / P(x))^2 dx \\
 &\leq \sum_{j \neq i} \sup_{\mathcal{D}_j} \int h_i(x) / h_j(x) P(x) dx = \sum_{j \neq i} e_{ij}(\cdot) \\
 &= 0.
 \end{aligned}$$

Summing up the results, we have

$$f(\cdot) \leq \max_{i,j} f_{ij}(\cdot) = 0.$$

Moreover, because

$$\lim_{\epsilon \rightarrow 0} \frac{f''(\cdot)}{\epsilon^p} = \lim_{\epsilon \rightarrow 0} \frac{f(\cdot)}{\epsilon^{\frac{p}{2}}} = 0;$$

we finally have  $f''(\cdot) = 0$  and thus  $f''(\cdot) = 0$ .

**Acknowledgments** \_\_\_\_\_

The work was supported by the HK RGC Earmarked Grant 4297/98E. We thank Professor Jiong Ruan and Dr. Yong Gao for the valuable discussions.

**References** \_\_\_\_\_

Delyon, B., Lavielle, M., & Moulines E. (1999). Convergence of a stochastic approximation version of the EM algorithm. *Annals of Statistics*, 27, 94–128.

- Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *J. R. Statist. Soc., B*, 39, 1–38.
- Gerald, S. R. (1980). *Matrix derivatives*. New York: Dekker.
- Jamshidian, M., & Jennrich R. I. (1997). Acceleration of the EM algorithm by using quasi-Newton methods. *J. R. Statist. Soc., B*, 59, 569–587.
- Jordan, M. I., & Jacobs, R. A. (1994). Hierarchical mixtures of experts and the EM algorithm. *Neural Computation*, 6, 181–214.
- Jordan, M. I., & Xu, L. (1995). Convergence results for the EP approach to mixtures of experts architectures. *Neural Networks*, 8(9), 1409–1431.
- Lange, K. (1995a). A gradient algorithm locally equivalent to the EM algorithm. *J. R. Statist. Soc., B*, 57, 425–437.
- Lange, K. (1995b). A quasi-Newton acceleration of the EM algorithm. *Statist. Sinica*, 5.
- Liu, C., & Rubin, D. B. (1994). The ECME algorithm: A simple extension of EM and ECM with faster monotone convergence. *Biometrika*, 81, 633–648.
- Meng, X. L. (1994). On the rate of convergence of the ECM algorithm. *Ann. Statist.*, 22, 326–339.
- Meng, X. L., & van Dyk D. (1997). The EM algorithm—An old folk-song sung to a fast new tune. *J. R. Statist. Soc., B*, 59, 511–567.
- Meng, X. L., & van Dyk D. (1998). Fast EM-type implementations for mixed effects models. *J. R. Statist. Soc., B*, 60, 559–578.
- Rabiner, L. R. (1989). A tutorial on hidden Markov models and selected applications in speech recognition. *Proc. of IEEE*, 77, 257–286.
- Redner, R. A., & Walker, H. F. (1984). Mixture densities, maximum likelihood, and the EM algorithm. *SIAM Review*, 26, 195–239.
- Veaux, D. De R. (1986). Parameter estimation for a mixture of linear regressions. Unpublished doctoral dissertation, Stanford University.
- Wu, C. F. J. (1983). On the convergence properties of the EM algorithm. *Annals of the Statistics*, 11, 95–103.
- Xu, L. (1997). Comparative analysis