

## Contents

1. Introduction	3
2. A normal mixture model with missing information	5
3. An efficient EM procedure for ML estimation	8
4. A data augmentation scheme for Bayesian sampling	10
5. Experimental results	15
6. Conclusions	21

## List of Tables

1	A comparison of CPU timings (in seconds) and relative reduced times (RRT) between GJ-EM algorithm (old) and our proposed procedure (new) under various missing rates. (Replications=500) . . . . .	16
2	A comparison of prediction accuracies for MI, EM and DA imputations with the standard deviations in parentheses for the <i>iris</i> data set. (Replications=500) . . . . .	17
3	A comparison of prediction accuracies for MI, EM and DA imputations with the standard deviations in parentheses for the <i>crabs</i> data set. (Replications=500) . . . . .	18
4	A comparison of average misclassification rates (%) between ML and Bayesian classifiers. (replicates=500) . . . . .	19

## List of Figures

1	ML and Bayesian density estimation for the two-component salmon data ( $\bullet$ , both attributes are completely observed; $\triangle$ , one of the two attributes is missing). . . . .	20
---	--	----

# On fast supervised learning for normal mixture models with missing information

Student: Hsiu J. Ho      Advisor: Dr. Tsung I. Lin

Department of Statistics  
Tunghai University  
Taichung  
Taiwan

## Abstract

It is an important research issue to deal with mixture models when missing values occur in the data. In this paper, computational strategies using auxiliary indicator matrices are introduced for handling mixtures of multivariate normal distributions in a more efficient manner, assuming that patterns of missingness are arbitrary and missing at random. We develop a novel structured EM algorithm which can dramatically save computation time and be exploited in many applications, such as density estimation, supervised clustering and prediction of missing values. In the aspect of multiple imputations for missing data, we also offer a data augmentation scheme using the Gibbs sampler. Our proposed methodologies are illustrated through some real data sets with varying proportions of missing values.

*Key words:* Bayesian classifier; Data augmentation; EM algorithm; Incomplete features; Rao-Blackwellization.

## 1. Introduction

Finite mixture models are known as powerful and flexible tools, which have been fully developed and applied in various theoretic and real problems as they are capable of modelling a wide range of densities, see monographs by Titterton et al. (1985), McLachlan and Basford (1988) and McLachlan and Peel (2000). However, missing values frequently appear in many real-world multivariate data sets that complicate data analyses and statistical inferences for practitioners. Missing data imputation techniques under the assumption of multivariate normal model have been well studied by Schafer (1997) and Liu (1999). In this decade, learning mixture models from incomplete data becomes an important research issue in multivariate analysis. The work on the use of Gaussian component was pioneered by Ghahramani and Jordan (1994), denoted by GJ hereafter. They present how to implement the Expectation- Maximization (EM) algorithm (Dempster et al. 1977) to compute maximum likelihood (ML) estimates from multivariate data with arbitrary pattern of missingness. They also compare the performance of EM imputation with a common mean imputation (MI) heuristic for the supervised classification of incomplete features.

Due to rapid advance of computational developments, Bayesian sampling-based approaches are usually considered as an alternative way in dealing with mixture models. There are plenty of papers in the literature to address the problem of fitting normal mixture models under Bayesian treatments. For example, Diebolt and Robert (1994) employ the data augmentation (DA) technique of Tanner and Wong (1987) as an approximation method for evaluating the posterior distribution

and showed a duality principle. Escobar and West (1995) present a nonparametric Bayesian density estimation for Dirichlet process mixture models. Richardson and Green (1997) and Zhang et al. (2004) propose a full Bayesian inference for a normal mixture model with unknown number of components using the reversible jump MCMC algorithm proposed by Green (1995). Stephens (2000) and Fruhwirth-Schnatter (2001) demonstrate Bayesian strategies for the elimination of label switching problems.

In this paper, we offer an efficient EM algorithm for the fitting of a likelihood-based normal mixture model using partially observed data. To reduce computational burden during the EM iterations, we incorporate two types of auxiliary binary indicator matrices corresponding to the observed and unobserved components of each datum. With strategies similar to EM, we also offer a DA computational technique for efficiently imputing missing values and learning parameters using the Gibbs sampler, which constructs a Markov chain that converges to a tractable posterior distribution (Geman and Geman, 1984). The feature of the chosen prior distributions are weakly informative to avoid mathematical and computational pitfalls of using improper priors in mixture model (Celex et al., 2000).

The rest of the paper proceeds as follows. In the next section, we describe the model and its notations, and present some important statistical properties based on the missing information framework. In Sections 3 and 4, two efficient EM and DA algorithms are developed to cope with ML and Bayesian estimation, respectively. We also investigate two issues regarding classification and prediction of incomplete features from ML and Bayesian perspectives. In Section 5, some real data sets are utilized to illustrate our proposed methodologies with varying proportions of

artificially missing values. Also, empirical comparisons between ML and Bayesian approaches in terms of classification and prediction accuracies for incomplete features are demonstrated. Finally, some concluding remarks are given in Section 6.

## 2. A normal mixture model with missing information

In the normal mixture model, we assume that  $\mathbf{Y} = (\mathbf{Y}_1, \dots, \mathbf{Y}_n)$  form a  $p$ -dimensional random sample from a population with  $g$  subclasses  $\mathcal{C}_1, \dots, \mathcal{C}_g$ , and each  $\mathbf{Y}_j$  has the density

$$f(\mathbf{Y}_j | \Theta) = \sum_{i=1}^g w_i \phi_p(\mathbf{Y}_j | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i), \quad w_i \geq 0, \quad \sum_{i=1}^g w_i = 1, \quad (1)$$

where  $w_i$ 's are mixing probabilities,  $\phi_p(\cdot | \boldsymbol{\mu}, \boldsymbol{\Sigma})$  denotes a  $p$ -dimensional multivariate normal component density with mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ , and  $\Theta = (w_1, \dots, w_g, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_g, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_g)$  is the vector of mixture model parameters subject to  $\sum_{i=1}^g w_i = 1$  and  $\boldsymbol{\Sigma}_i$ 's are positive definite matrices. Thus, there are  $g(p+1)(p+2)/2 - 1$  distinct parameters in model (1).

Typically, in the EM framework, mixture models can be characterized as having an incomplete data structure. It is convenient to formalize the missing part as a set of membership labels  $\mathbf{Z} = (\mathbf{Z}_1, \dots, \mathbf{Z}_n)$  with each label  $\mathbf{Z}_j = (Z_{1j}, \dots, Z_{gj})$  being a binary vector such that  $Z_{ij} = 1$  if  $\mathbf{Y}_j$  belongs to component  $i$  and  $Z_{ij} = 0$  otherwise. Given the mixing probabilities  $\boldsymbol{\omega}$ ,  $\mathbf{Z}_1, \dots, \mathbf{Z}_n$  independently follow a multinomial distribution. We shall write  $\mathbf{Z}_j \sim \mathcal{M}(1; w_1, \dots, w_g)$ .

For notational simplicity, let

$$\Delta_{ij} = (\mathbf{Y}_j - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}_i^{-1} (\mathbf{Y}_j - \boldsymbol{\mu}_i), \quad (2)$$

denote the Mahalanobis distance for  $\mathbf{Y}_j$  with respect to mean  $\boldsymbol{\mu}_i$  and covariance

matrix  $\Sigma_i$ . The complete likelihood function for  $\Theta$  is

$$L_c(\Theta|\mathbf{Y}, \mathbf{Z}) \propto \prod_{j=1}^n \prod_{i=1}^g \left( w_i |\Sigma_i|^{-1/2} \exp\left(-\frac{1}{2} \Delta_{ij}\right) \right)^{Z_{ij}}. \quad (3)$$

We consider the maximum likelihood estimation problem of model (1) when  $\mathbf{Y}$  are not completely observed. We further assume that the patterns of missingness are arbitrary and missing at random (MAR), see Rubin (1976) and Little and Rubin (2002) for more details. Generally speaking, MAR refers to the missingness depends only on observed values but not on missing values.

Let  $\mathbf{Y}_j$  be partitioned into two components  $(\mathbf{Y}_j^o, \mathbf{Y}_j^m)$ , where  $\mathbf{Y}_j^o$  ( $p_j^o \times 1$ ) and  $\mathbf{Y}_j^m$  ( $(p - p_j^o) \times 1$ ) denote the observed and missing components of  $\mathbf{Y}_j$ , respectively. To facilitate the EM algorithm, it is advantageous to introduce two types of binary indicator matrices, denoted by  $\mathbf{O}_j$  and  $\mathbf{M}_j$  hereafter, corresponding to  $\mathbf{Y}_j$  such that  $\mathbf{Y}_j^o = \mathbf{O}_j \mathbf{Y}_j$  and  $\mathbf{Y}_j^m = \mathbf{M}_j \mathbf{Y}_j$ , respectively. Notice that  $\mathbf{O}_j$  and  $\mathbf{M}_j$  are  $p_j^o \times p$  and  $(p - p_j^o) \times p$  matrices extracted from a  $p$ -dimensional identity matrix  $\mathbf{I}_p$  corresponding to row-positions of  $\mathbf{Y}_j^o$  and  $\mathbf{Y}_j^m$  in  $\mathbf{Y}_j$ , respectively. We then have the following propositions.

**Proposition 1.** *Suppose  $\mathbf{Y}_j$  is partitioned into two components  $(\mathbf{Y}_j^o, \mathbf{Y}_j^m)$ , where  $\mathbf{Y}_j^o = \mathbf{O}_j \mathbf{Y}_j$  and  $\mathbf{Y}_j^m = \mathbf{M}_j \mathbf{Y}_j$ . We thus have*

$$\mathbf{Y}_j = \begin{cases} \mathbf{Y}_j^o, & \text{if } p_j^o = p; \\ \mathbf{O}_j^\top \mathbf{Y}_j^o + \mathbf{M}_j^\top \mathbf{Y}_j^m, & \text{if } 1 \leq p_j^o < p, \end{cases}$$

and  $\mathbf{O}_j^\top \mathbf{O}_j + \mathbf{M}_j^\top \mathbf{M}_j = \mathbf{I}_p$ .

**Proof:** The proof is straightforward and hence is omitted.

**Proposition 2.** Let  $\mathbf{Y}_j \sim \sum_{i=1}^g w_i \phi_p(\mathbf{Y}_j | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ , and let  $\mathbf{Y}_j^o$  and  $\mathbf{Y}_j^m$  be the observed and missing components corresponding  $\mathbf{Y}_j$ , respectively. The marginal distribution of  $\mathbf{Y}_j^o$  is denoted by  $\mathbf{Y}_j^o \sim \sum_{i=1}^g w_i \phi_{p_j^o}(\mathbf{Y}_j^o | \boldsymbol{\mu}_{ij}^o, \boldsymbol{\Sigma}_{ij}^{oo})$ , where

$$\phi_{p_j^o}(\mathbf{Y}_j^o | \boldsymbol{\mu}_{ij}^o, \boldsymbol{\Sigma}_{ij}^{oo}) = (2\pi)^{-p_j^o/2} |\boldsymbol{\Sigma}_{ij}^{oo}|^{-1/2} \exp(-\frac{1}{2} \Delta_{ij}^o),$$

is the component density, and

$$\begin{aligned} \boldsymbol{\mu}_{ij}^o &= \mathbf{O}_j \boldsymbol{\mu}_i, \quad \boldsymbol{\Sigma}_{ij}^{oo} = \mathbf{O}_j \boldsymbol{\Sigma}_i \mathbf{O}_j^\top, \quad \Delta_{ij}^o = (\mathbf{Y}_j - \boldsymbol{\mu}_i)^\top \mathbf{S}_{ij}^{oo} (\mathbf{Y}_j - \boldsymbol{\mu}_i), \\ \mathbf{S}_{ij}^{oo} &= \mathbf{O}_j^\top (\mathbf{O}_j \boldsymbol{\Sigma}_i \mathbf{O}_j^\top)^{-1} \mathbf{O}_j. \end{aligned} \quad (4)$$

Consequently,  $\mathbf{Y}_j^m | \mathbf{Y}_j^o \sim \sum_{i=1}^g w_{ij}^* \phi_{p-p_j^o}(\mathbf{Y}_j^m | \boldsymbol{\mu}_{ij}^{m\cdot o}, \boldsymbol{\Sigma}_{ij}^{mm\cdot o})$ , where

$$\phi_{p-p_j^o}(\mathbf{Y}_j^m | \boldsymbol{\mu}_{ij}^{m\cdot o}, \boldsymbol{\Sigma}_{ij}^{mm\cdot o}) = (2\pi)^{-(p-p_j^o)/2} |\boldsymbol{\Sigma}_{ij}^{mm\cdot o}|^{-1/2} \exp(-\frac{1}{2} \Delta_{ij}^{m\cdot o}),$$

and

$$\begin{aligned} w_{ij}^* &= w_i \phi_{p_j^o}(\mathbf{Y}_j^o | \boldsymbol{\mu}_{ij}^o, \boldsymbol{\Sigma}_{ij}^{oo}) / \sum_{h=1}^g w_h \phi_{p_j^o}(\mathbf{Y}_j^o | \boldsymbol{\mu}_{hj}^o, \boldsymbol{\Sigma}_{hj}^{oo}), \\ \boldsymbol{\mu}_{ij}^{m\cdot o} &= \mathbf{M}_j (\boldsymbol{\mu}_i + \boldsymbol{\Sigma}_i \mathbf{S}_{ij}^{oo} (\mathbf{Y}_j - \boldsymbol{\mu}_i)), \quad \boldsymbol{\Sigma}_{ij}^{mm\cdot o} = \mathbf{E}_{ij} \boldsymbol{\Sigma}_i \mathbf{M}_j^\top, \\ \mathbf{E}_{ij} &= \mathbf{M}_j (\mathbf{I}_p - \boldsymbol{\Sigma}_i \mathbf{S}_{ij}^{oo}), \quad \Delta_{ij}^{m\cdot o} = (\mathbf{Y}_j - \boldsymbol{\mu}_i)^\top \mathbf{S}_{ij}^{mm\cdot o} (\mathbf{Y}_j - \boldsymbol{\mu}_i), \\ \mathbf{S}_{ij}^{mm\cdot o} &= \mathbf{E}_{ij}^\top (\mathbf{E}_{ij} \boldsymbol{\Sigma}_i \mathbf{M}_j^\top)^{-1} \mathbf{E}_{ij}. \end{aligned} \quad (5)$$

**Proof:** The sketch of the proof is given in Appendix A.

To enhance the computational efficiency for estimation, we suggest to rearrange  $\mathbf{Y}$  according to unique missing patterns of the data. The procedure can be implemented as follows:

- (a) Build a binary indicator matrix,  $\mathbf{R} = [r_{ij}]_{n \times p}$ , with each entry  $r_{ij} = 1$  if  $Y_{ij}$  is missing and  $r_{ij} = 0$  otherwise.

- (b) Let  $\mathbf{z} = \mathbf{R} \times \mathbf{b}$ , where  $\mathbf{b} = (2^1, 2^2, \dots, 2^p)^\top$ . Note that the number of unique missing patterns is equal to the number of unique elements in  $\mathbf{z}$ .
- (c) Denoting  $\mathbf{z}^*$  by sorting  $\mathbf{z}$  in an ascending or descending order, we then rearrange  $\mathbf{Y}$  according to the row positions of  $\mathbf{z}^*$  in  $\mathbf{z}$ . This will yield clustering of identical patterns of missingness in  $\mathbf{Y}$  which are adjacent to each other.

### 3. An efficient EM procedure for ML estimation

Let  $\mathbf{Y}^o = (\mathbf{Y}_1^o, \dots, \mathbf{Y}_n^o)$  and  $\mathbf{Y}^m = (\mathbf{Y}_1^m, \dots, \mathbf{Y}_n^m)$  denote the observed portion and missing portion of the data, respectively. The complete-data log-likelihood function can be reexpressed by

$$\begin{aligned}
& \ell_c(\boldsymbol{\Theta} | \mathbf{Y}^o, \mathbf{Y}^m, \mathbf{Z}) \\
&= \ell_{c_1}(\boldsymbol{\omega} | \mathbf{Y}^o, \mathbf{Y}^m, \mathbf{Z}) + \ell_{c_2}(\boldsymbol{\Psi} | \mathbf{Y}^o, \mathbf{Y}^m, \mathbf{Z}) \\
&= \sum_{i=1}^g \sum_{j=1}^n Z_{ij} \log w_i + \frac{1}{2} \sum_{i=1}^g \left( \log |\boldsymbol{\Sigma}^{-1}| \sum_{j=1}^n Z_{ij} - \sum_{j=1}^n Z_{ij} (\Delta_{ij}^o + \Delta_{ij}^{m \cdot o}) \right), \quad (6)
\end{aligned}$$

where  $\boldsymbol{\omega} = (w_1, \dots, w_g)$  and  $\boldsymbol{\Psi} = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_g, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_g)$ . From (5), it is easy to verify that  $\boldsymbol{\Sigma}_i^{-1} = \mathbf{S}_{ij}^{oo} + \mathbf{S}_{ij}^{mm \cdot o}$  and  $\mathbf{O}_j^\top \mathbf{O}_j (\mathbf{I}_p - \boldsymbol{\Sigma}_i \mathbf{S}_{ij}^{oo}) = \mathbf{0}$ . Hence, we have the following result.

**Proposition 3.** *The conditional expectation of (6) is give by*

$$\begin{aligned}
Q(\boldsymbol{\Theta} | \hat{\boldsymbol{\Theta}}^{(k)}) &= E(\ell_c(\boldsymbol{\Theta} | \mathbf{Y}^o, \mathbf{Y}^m, \mathbf{Z}) | \mathbf{Y}^o, \hat{\boldsymbol{\Theta}}^{(k)}) \\
&= Q_1(\boldsymbol{\omega} | \hat{\boldsymbol{\Theta}}^{(k)}) + Q_2(\boldsymbol{\Psi} | \hat{\boldsymbol{\Theta}}^{(k)}).
\end{aligned}$$



It follows that

$$Q_1(\mathbf{w}|\hat{\Theta}^{(k)}) = \sum_{i=1}^g \sum_{j=1}^n \hat{Z}_{ij}^{(k)} \log w_i, \quad (7)$$

$$Q_2(\Psi|\hat{\Theta}^{(k)}) = \frac{1}{2} \sum_{i=1}^g \left( \log |\Sigma_i^{-1}| \sum_{j=1}^n \hat{Z}_{ij}^{(k)} - \text{tr} \left( \Sigma_i^{-1} \sum_{j=1}^n \Omega_{ij}^{(k)} \right) \right), \quad (8)$$

where

$$\Omega_{ij}^{(k)} = \hat{Z}_{ij}^{(k)} \left( (\hat{\mathbf{Y}}_{ij}^{(k)} - \boldsymbol{\mu}_i)(\hat{\mathbf{Y}}_{ij}^{(k)} - \boldsymbol{\mu}_i)^\top + (\mathbf{I}_p - \hat{\Sigma}_i^{(k)} \hat{\mathbf{S}}_{ij}^{\text{oo}(k)}) \hat{\Sigma}_i^{(k)} \right), \quad (9)$$

$$\hat{Z}_{ij}^{(k)} = \frac{\hat{w}_i^{(k)} \phi_{p_j^o}(\mathbf{Y}_j^o | \hat{\boldsymbol{\mu}}_{ij}^{\text{o}(k)}, \hat{\Sigma}_{ij}^{\text{oo}(k)})}{\sum_{h=1}^g \hat{w}_h^{(k)} \phi_{p_j^o}(\mathbf{Y}_j^o | \hat{\boldsymbol{\mu}}_{hj}^{\text{o}(k)}, \hat{\Sigma}_{hj}^{\text{oo}(k)})}, \quad (10)$$

$$\hat{\mathbf{Y}}_{ij}^{(k)} = \hat{\boldsymbol{\mu}}_i^{(k)} + \hat{\Sigma}_i^{(k)} \hat{\mathbf{S}}_{ij}^{\text{oo}(k)} (\mathbf{Y}_j - \hat{\boldsymbol{\mu}}_i^{(k)}), \quad (11)$$

and  $\hat{\mathbf{S}}_{ij}^{\text{oo}(k)}$  is  $\mathbf{S}_{ij}^{\text{oo}}$  given in (4) with  $\Sigma_i$  replaced by  $\hat{\Sigma}_i^{(k)}$ .

**Proof:** The detailed proof is shown in Appendix B.

By these propositions, a modified version of GJ's EM algorithm can be implemented as follows:

**E-step:** Given  $\Theta = \hat{\Theta}^{(k)}$ , impute  $\hat{Z}_{ij}^{(k)}$  and  $\hat{\mathbf{Y}}_{ij}^{(k)}$  for  $i = 1, \dots, g$  and  $j = 1, \dots, n$ , using (10) and (11).

**M-Step:**

1. Update  $\hat{w}_i^{(k+1)}$  by maximizing (7) over  $w_i$ , which leads to

$$\hat{w}_i^{(k+1)} = \frac{1}{n} \sum_{j=1}^n \hat{Z}_{ij}^{(k)}.$$

2. Fix  $\Sigma_i$  at  $\hat{\Sigma}_i^{(k)}$ , update  $\hat{\boldsymbol{\mu}}_i^{(k+1)}$  by maximizing (8) over  $\boldsymbol{\mu}_i$ , which leads to

$$\hat{\boldsymbol{\mu}}_i^{(k+1)} = \frac{\sum_{j=1}^n \hat{Z}_{ij}^{(k)} \hat{\mathbf{Y}}_{ij}^{(k)}}{\sum_{j=1}^n \hat{Z}_{ij}^{(k)}}.$$

3. Fix  $\boldsymbol{\mu}_i$  at  $\hat{\boldsymbol{\mu}}_i^{(k+1)}$ , update  $\hat{\boldsymbol{\Sigma}}_i^{(k+1)}$  by maximizing constrained (8) over  $\boldsymbol{\Sigma}_i$ , which leads to

$$\hat{\boldsymbol{\Sigma}}_i^{(k+1)} = \frac{\sum_{j=1}^n \hat{\boldsymbol{\Omega}}_{ij}^{(k)}}{\sum_{j=1}^n \hat{Z}_{ij}^{(k)}},$$

where  $\hat{\boldsymbol{\Omega}}_{ij}^{(k)}$  is  $\boldsymbol{\Omega}_{ij}^{(k)}$  in (9) with  $\boldsymbol{\mu}_i$  replaced by  $\hat{\boldsymbol{\mu}}_i^{(k+1)}$ .

We remark two major advantages of the above EM algorithm:

- (a) With auxiliary matrices  $\boldsymbol{O}_j$ 's obtained at the initiation, there is no need to take care of the associated row positions of missing values at each iteration.
- (b) The implementation of M-step has low computational cost as it is similar to the case of no missing values. Therefore, the modified EM algorithm is more straightforward than the version of GJ.

Applying Bayes' theorem, the posterior probability of the  $\mathbf{Y}_j$  belonging to  $\mathcal{C}_i$  can be estimated by

$$\hat{w}_{ij}^* = \Pr(Z_{ij} = 1 | \mathbf{Y}^o, \hat{\boldsymbol{\Theta}}) = \frac{\hat{w}_i \phi_{p_j^o}(\mathbf{Y}_j^o | \hat{\boldsymbol{\mu}}_{ij}^o, \hat{\boldsymbol{\Sigma}}_{ij}^{oo})}{\sum_{h=1}^g \hat{w}_h \phi_{p_j^o}(\mathbf{Y}_j^o | \hat{\boldsymbol{\mu}}_{hj}^o, \hat{\boldsymbol{\Sigma}}_{hj}^{oo})}. \quad (12)$$

By the ML classification theory (Basford and McLachlan, 1985),  $\mathbf{Y}_j$  is assigned to  $\mathcal{C}_s$  if  $\hat{w}_{sj}^* > \hat{w}_{ij}^*$  ( $i = 1, \dots, g; i \neq s$ ).

Consequently, an ML predictor for the missing component  $\mathbf{Y}_j^m$  is given by

$$\hat{\mathbf{Y}}_j^m = E(\mathbf{Y}_j^m | \mathbf{Y}^o, \hat{\boldsymbol{\Theta}}) = \mathbf{M}_j \sum_{i=1}^g \hat{w}_{ij}^* \left( \hat{\boldsymbol{\mu}}_i + \hat{\boldsymbol{\Sigma}}_i \hat{\mathbf{S}}_{ij}^{oo} (\mathbf{Y}_j - \hat{\boldsymbol{\mu}}_i) \right). \quad (13)$$

#### 4. A data augmentation scheme for Bayesian sampling

The Data augmentation (DA) of Tanner and Wong (1987) is a general and effective algorithm for producing multiple imputation of missing data. The DA has been

broadly applied in a variety of missing data problems, see for example, Schafer (1997), Little and Rubin (2002) and references therein. In this section, we construct an efficient DA algorithm that combines the latent variables  $\mathbf{Z}$  and unobserved data  $\mathbf{Y}^m$  for simulating the posterior density of  $\Theta$ .

The DA algorithm consists of the imputation step (I-step) and the posterior step (P-step). At the  $k$ th iteration, the I-step is defined by drawing imputations of  $\mathbf{Z}_j$ 's and  $\mathbf{Y}_j^m$ 's from the predictive distributions  $p(\mathbf{Z}_j | \mathbf{Y}^o, \Theta^{(k)})$  and  $p(\mathbf{Y}_j^m | \mathbf{Y}^o, \mathbf{Z}_j, \Theta^{(k)})$ , respectively, and the P-step refer to generating  $\Theta^{(k+1)}$  from  $p(\Theta | \mathbf{Y}^o, \mathbf{Y}^{m^{(k+1)}}, \mathbf{Z}^{(k+1)})$ . To perform the Bayesian inference for mixture models, it is necessary to choose a proper prior distribution for each parameter to avoid yielding improper posterior distributions (Celeux et al., 2000). In univariate normal mixture models, Diebolt and Robert (1994) and Richardson and Green (1997) have suggested some conjugate prior distributions. In the context of multivariate version with missing information, our chosen priors for the model parameters primarily follow the suggestion of Stephens (2000). They are given by

$$\begin{aligned} \mathbf{w} &\sim \mathcal{D}(\delta, \dots, \delta), \\ \boldsymbol{\mu}_i &\sim \mathcal{N}_p(\boldsymbol{\xi}, \boldsymbol{\kappa}^{-1}) \quad (i = 1, \dots, g), \\ \boldsymbol{\Sigma}_i^{-1} | \mathbf{B} &\sim \mathcal{W}_p(2\alpha, (2\mathbf{B})^{-1}) \quad (i = 1, \dots, g), \\ \mathbf{B} &\sim \mathcal{W}_p(2\gamma, (2\mathbf{H})^{-1}), \end{aligned}$$

where  $\mathbf{B}$  is a  $p \times p$  hyperparameter matrix,  $\boldsymbol{\kappa}$  and  $\mathbf{H}$  are  $p \times p$  constant matrices,  $\boldsymbol{\xi}$  is a  $p \times 1$  constant vector,  $\alpha$ ,  $\delta$  and  $\gamma$  are constant scalars,  $\mathcal{D}(\delta, \dots, \delta)$  denotes the symmetric Dirichlete distribution with density

$$f(\mathbf{w}|\delta) = \frac{\Gamma(g\delta)}{\Gamma(\delta)^g} w_1^{\delta-1} \dots w_{g-1}^{\delta-1} (1 - w_1 - \dots - w_{g-1})^{\delta-1},$$

and  $\mathcal{W}_p(\nu, \mathbf{A})$  denotes the Wishart distribution with density

$$f(\mathbf{U}|\nu, \mathbf{A}) \propto |\mathbf{A}|^{-\nu/2} |\mathbf{U}|^{(\nu-p-1)/2} \exp\left(-\frac{1}{2}\text{tr}(\mathbf{U}\mathbf{A}^{-1})\right).$$

The joint prior distribution function of  $\Theta$  and  $\mathbf{B}$  is

$$\begin{aligned} \pi(\Theta, \mathbf{B}) &\propto w_1^{\delta-1} \dots w_g^{\delta-1} |\mathbf{B}|^{g\alpha+(2\gamma-p-1)/2} \exp\{-\text{tr}(\mathbf{H}\mathbf{B})\} \\ &\times \prod_{i=1}^g |\Sigma_i^{-1}|^{(2\alpha-p-1)/2} \exp\left(-\frac{1}{2}(\boldsymbol{\mu}_i - \boldsymbol{\xi})^\top \boldsymbol{\kappa}(\boldsymbol{\mu}_i - \boldsymbol{\xi}) - \text{tr}(\mathbf{B}\Sigma_i^{-1})\right). \end{aligned} \quad (14)$$

Upon multiplying (3) and (14), we have the following joint posterior density:

$$\begin{aligned} p(\Theta, \mathbf{B}, \mathbf{Y}^m, \mathbf{Z}|\mathbf{Y}^o) &\propto w_1^{\delta-1} \dots w_g^{\delta-1} |\mathbf{B}|^{g\alpha+(2\gamma-p-1)/2} \exp(-\text{tr}(\mathbf{H}\mathbf{B})) \\ &\times \prod_{i=1}^g \exp\left(-\frac{1}{2}(\boldsymbol{\mu}_i - \boldsymbol{\xi})^\top \boldsymbol{\kappa}(\boldsymbol{\mu}_i - \boldsymbol{\xi})\right) |\Sigma_i^{-1}|^{(2\alpha-p-1)/2} \exp(-\text{tr}(\mathbf{B}\Sigma_i^{-1})) \\ &\times \prod_{j=1}^n \prod_{i=1}^g \left(w_i |\Sigma_i^{-1}|^{1/2} \exp\left(-\frac{1}{2}(\Delta_{ij}^o + \Delta_{ij}^{m\cdot o})\right)\right)^{Z_{ij}}, \end{aligned} \quad (15)$$

where  $\Delta_{ij}^o$  and  $\Delta_{ij}^{m\cdot o}$  are given in (4) and (5), respectively.

**Proposition 4.** *The full conditional posteriors of  $\Theta$ ,  $\mathbf{B}$ ,  $\mathbf{Z}$  and  $\mathbf{Y}^m$  are as follows (the symbol “ $|\dots$ ” denotes conditioning on all other variables):*

$$\begin{aligned} p(\mathbf{Z}_j|\mathbf{Y}^o, \Theta) &\propto \prod_{i=1}^g \left(w_i \phi_{p_j^o}(\mathbf{Y}_j^o | \boldsymbol{\mu}_{ij}^o, \Sigma_{ij}^{oo})\right)^{Z_{ij}}, \\ p(\mathbf{Y}_j^m | Z_{ij} = 1, \dots) &\propto \exp\left(-\frac{1}{2}(\mathbf{Y}_j^m - \boldsymbol{\mu}_{ij}^{m\cdot o})^\top \Sigma_{ij}^{mm\cdot o^{-1}} (\mathbf{Y}_j^m - \boldsymbol{\mu}_{ij}^{m\cdot o})\right), \\ p(\mathbf{w}|\dots) &\propto \prod_{i=1}^g w_i^{\sum_{j=1}^n Z_{ij} + \delta - 1}, \\ p(\boldsymbol{\mu}_i|\dots) &\propto \exp\left(-\frac{1}{2}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_i^*)^\top \Sigma_i^{*-1} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_i^*)\right), \\ p(\mathbf{B}|\dots) &\propto |\mathbf{B}|^{(2(g\alpha+\gamma)-p-1)/2} \exp\left(-\text{tr}\left(\mathbf{B}\left(\mathbf{H} + \sum_{i=1}^g \Sigma_i^{-1}\right)\right)\right), \\ p(\Sigma_i^{-1}|\dots) &\propto |\Sigma_i^{-1}|^{(\alpha^*-p-1)/2} \exp\left(-\frac{1}{2}\text{tr}(\Sigma_i^{-1}\mathbf{A}_i)\right), \end{aligned}$$

where  $\boldsymbol{\mu}_{ij}^{\text{m}\cdot\text{o}}$  and  $\boldsymbol{\Sigma}_{ij}^{\text{mm}\cdot\text{o}}$  are given by (5), and

$$\boldsymbol{\Sigma}_i^* = \left( \boldsymbol{\Sigma}_i^{-1} \sum_{j=1}^n Z_{ij} + \boldsymbol{\kappa} \right)^{-1}, \quad (16)$$

$$\boldsymbol{\mu}_i^* = \boldsymbol{\Sigma}_i^* \left( \boldsymbol{\Sigma}_i^{-1} \sum_{j=1}^n Z_{ij} \mathbf{Y}_j + \boldsymbol{\kappa} \boldsymbol{\xi} \right), \quad (17)$$

$$\alpha_i^* = \sum_{j=1}^n Z_{ij} + 2\alpha, \quad (18)$$

$$\mathbf{A}_i = 2\mathbf{B} + \sum_{j=1}^n Z_{ij} (\mathbf{Y}_j - \boldsymbol{\mu}_i) (\mathbf{Y}_j - \boldsymbol{\mu}_i)^\top, \quad (19)$$

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

**Proof:** The proof is straightforward and hence is omitted.

In the simulation process, samples for  $\mathbf{Z}$ ,  $\mathbf{Y}^{\text{m}}$ ,  $\mathbf{B}$  and  $\boldsymbol{\Theta}$  are alternately generated, the DA algorithm using the Gibbs sampler can be implemented as follows:

**I-Step:**

1. Given  $\boldsymbol{\Theta}$ ,  $\mathbf{Y}^{\text{m}}$  and  $\mathbf{Y}^{\text{o}}$ , generate  $\mathbf{Z}_j$  from  $\mathcal{M}(1; r_{1j}, \dots, r_{gj})$ , where

$$r_{ij} = \frac{w_i \phi_{p_j^{\text{o}}}(\mathbf{Y}_j^{\text{o}} | \boldsymbol{\mu}_{ij}^{\text{o}}, \boldsymbol{\Sigma}_{ij}^{\text{oo}})}{\sum_{s=1}^g w_s \phi_{p_j^{\text{o}}}(\mathbf{Y}_j^{\text{o}} | \boldsymbol{\mu}_{sj}^{\text{o}}, \boldsymbol{\Sigma}_{sj}^{\text{oo}})}.$$

2. Generate  $\mathbf{Y}_j^{\text{m}}$  given  $Z_{ij} = 1$ ,  $\boldsymbol{\Theta}$  and  $\mathbf{Y}^{\text{o}}$ , from  $N_{p-p_j^{\text{o}}}(\boldsymbol{\mu}_{ij}^{\text{m}\cdot\text{o}}, \boldsymbol{\Sigma}_{ij}^{\text{mm}\cdot\text{o}})$ , where  $\boldsymbol{\mu}_{ij}^{\text{m}\cdot\text{o}}$  and  $\boldsymbol{\Sigma}_{ij}^{\text{mm}\cdot\text{o}}$  are as in (5).

**P-Step:**

1. Generate  $\mathbf{w}$  given  $\mathbf{Z}$  from  $\mathcal{D}(n_1 + \delta, \dots, n_g + \delta)$ , where  $n_i = \sum_{j=1}^n Z_{ij}$ .
2. Generate  $\boldsymbol{\mu}_i$  given  $\mathbf{Z}$ ,  $\boldsymbol{\Sigma}_i$ ,  $\mathbf{Y}^{\text{o}}$  and  $\mathbf{Y}^{\text{m}}$  from  $\mathcal{N}_p(\boldsymbol{\mu}_i^*, \boldsymbol{\Sigma}_i^*)$  with  $\boldsymbol{\mu}_i^*$  and  $\boldsymbol{\Sigma}_i^*$  given in (17) and (16), respectively.

3. Generate  $\mathbf{B}$  given  $\Sigma_1, \dots, \Sigma_g$  from  $\mathcal{W}_p(2\gamma^*, (2\mathbf{H}^*)^{-1})$ , where  $\gamma^* = g\alpha + \gamma$  and

$$\mathbf{H}^* = \mathbf{H} + \sum_{i=1}^g \Sigma_i^{-1}.$$

4. Generate  $\Sigma_i^{-1}$  given  $\mathbf{Z}, \boldsymbol{\mu}_i, \mathbf{Y}^o$  and  $\mathbf{Y}^m$  from  $\mathcal{W}_p(\alpha_i^*, \mathbf{A}_i^{-1})$ , where  $\alpha_i^*$  and  $\mathbf{A}_i$  are given in (18) and (19), respectively.

To satisfy the ‘‘Principle of Stable Estimation’’ of Edwards et al. (1963) in the Bayesian treatment, we need to specify  $(\boldsymbol{\xi}, \boldsymbol{\kappa}, \alpha, \gamma, \mathbf{H})$  so as to be insensitive to changes of the prior. Specifically, it is often to choose  $\delta = 1$ . For  $\boldsymbol{\xi}$  and  $\boldsymbol{\kappa}$ , we let  $\boldsymbol{\xi}$  be the empirical mean vector and  $\boldsymbol{\kappa}^{-1} = (1 - \eta)^{-1} \text{diag}\{R_1^2, \dots, R_p^2\}$ , where  $\eta$  is the percentage of missing values of the data which is used to adjust the flatness and  $R_i$  is the range of the observed values of variable  $i$ . As a generalization of Richard and Green (1997), we take  $\alpha = p + 1$ ,  $\gamma = (p + 1)/10$  and  $\mathbf{H} = 10\boldsymbol{\kappa}$ .

We are interested in the classification and prediction problems for incomplete features. Under certain conditions, quantites based on Rao-Blackwellization (Gelfand and Smith, 1990) often greatly improve the precision of Monte Carlo estimates. Given a set of converged Monte Carlo DA samples  $\Theta^{(\ell)}$  ( $\ell = 1, \dots, L$ ), a Bayesian predictor for  $\mathbf{Y}_j^m$  is given by

$$\begin{aligned} \tilde{\mathbf{Y}}_j^m &= \frac{1}{L} \sum_{\ell=1}^L \mathbb{E}(\mathbf{Y}_j^m | \mathbf{Y}_j^o, \Theta^{(\ell)}) \\ &= \mathbf{M}_j \frac{1}{L} \sum_{\ell=1}^L \left( \sum_{i=1}^g r_{ij}^{(\ell)} \left( \boldsymbol{\mu}_i^{(\ell)} + \Sigma_i^{(\ell)} \mathbf{S}_{ij}^{\text{oo}(\ell)} (\mathbf{Y}_j - \boldsymbol{\mu}_i^{(\ell)}) \right) \right), \end{aligned} \quad (20)$$

where

$$r_{ij}^{(\ell)} = \frac{w_i^{(\ell)} \phi_{p_j^o}(\mathbf{Y}_j^o | \hat{\boldsymbol{\mu}}_i^{o(\ell)}, \Sigma_i^{\text{oo}(\ell)})}{\sum_{h=1}^g w_h^{(\ell)} \phi_{p_j^o}(\mathbf{Y}_j^o | \boldsymbol{\mu}_h^{o(\ell)}, \Sigma_h^{\text{oo}(\ell)})}.$$

Consequently, a Bayesian classifier for  $\mathbf{Y}_j$  can be estimated by averaging over

the draws of  $\Theta^{(\ell)}$

$$\hat{r}_{ij}^* = \Pr(Z_{ij} = 1 | \mathbf{Y}_j^o) \approx \frac{1}{L} \sum_{\ell=1}^L r_{ij}^{(\ell)}. \quad (21)$$

By the Bayesian classification rule,  $\mathbf{Y}_j$  is assigned to  $\mathcal{C}_s$  if  $\hat{r}_{sj}^* > \hat{r}_{ij}^*$  ( $i = 1, \dots, g; i \neq s$ ).

## 5. Experimental results

For illustration purposes, we start to apply results developed in Sections 2-4 to two famous multivariate data sets. One is the *iris* data taken from Anderson (1935) or Fisher (1936). It consists of 4-dimensional measurements in centimeters on the attributes of petal length, petal width, sepal length and sepal width for 50 flower specimens of each of three species: *setosa*, *versicolor*, and *virginica*. The other is the *crabs* data of Campbell and Mahon (1974) on the genus *Leptograpsus*. It consists of 5-dimensional morphological measurements on the attributes of width of frontal lip, rear width, length along the mid-line of the carapace, maximum width of the carapace and body depth for 50 crabs of each of four groups: blue male, blue female, orange male and orange female. Both data sets are included as a part of the R package, which is freely available at the web site <http://cran.r-project.org>.

To conduct experimental studies, we first generate 500 artificially missing data sets by deleting at random from the three data sets under various specified missing rate  $\eta$  (proportion of missing values) while we maintain each datum that has at least one observed attribute. Table 1 presents the computation times of our developed EM algorithm and those of using GJ-EM. All computations are solely carried out by R package in the environment of a desktop PC (CPU: 3G-MHz/Intel Pentium 4 Processor; RAM: 1024 MB/DDR-400). Since the programming implementations

Table 1: A comparison of CPU timings (in seconds) and relative reduced times (RRT) between GJ-EM algorithm (old) and our proposed procedure (new) under various missing rates. (Replications=500)

Data	$\eta = 10\%$			$\eta = 20\%$			$\eta = 30\%$		
	old	new	RRT	old	new	RRT	old	new	RRT
<i>iris</i>	12.47	1.22	90.2%	21.51	1.61	92.5%	56.21	3.61	93.6%
<i>crabs</i>	34.72	3.27	90.6%	78.77	6.78	91.4%	265.01	20.68	92.2%

$$\text{RRT} = (\text{old} - \text{new}) / \text{old} \times 100\%$$

have many characteristics (e.g., vector or matrix subroutines instead of loops), the CPU times in Table 1 might not be directly comparable, but provide a sense of their actual performances in a practical setting. As seen in the table, all computation times are dramatically reduced over 90% by using the new EM procedure.

To exemplify the predictive performance for the EM and DA imputation methods, see Equations (13) and (20), together with the traditional mean imputation (MI) method, known as “filling-in” with the sample mean of the associated attribute, we utilize the pseudo-cross-validation (PSV) of Stone (1974) to evaluate these three approaches. A relative tolerance of  $10^{-8}$  for the log-likelihood function and parameter estimates are used as the convergence criterion for the EM algorithm. As for the DA algorithm, we take the ML estimates as the initialization and carry out 2,000 iterations with the first 1,000 iterations as burn-in and the remaining 1,000 iterations as inference samples. It is noted that our chosen burn-in number is much larger than needed based on checking the *multivariate potential scale reduction factor* (MPSRF) of Brooks and Gelman (1998). As for discrepancy measures, we use the mean absolute error (MAE), the mean absolute relative error (MARE) and root



Table 2: A comparison of prediction accuracies for MI, EM and DA imputations with the standard deviations in parentheses for the *iris* data set. (Replications=500)

$\eta$	MAE			MARE			RMSE		
	MI	EM	DA	MI	EM	DA	MI	EM	DA
10%	0.812 (0.081)	0.213 (0.026)	0.210 (0.026)	0.697 (0.186)	0.100 (0.027)	0.099 (0.027)	1.062 (0.096)	0.285 (0.050)	0.280 (0.050)
20%	0.816 (0.053)	0.237 (0.025)	0.233 (0.025)	0.675 (0.129)	0.114 (0.031)	0.113 (0.031)	1.071 (0.065)	0.331 (0.060)	0.326 (0.060)
30%	0.820 (0.046)	0.268 (0.023)	0.259 (0.022)	0.684 (0.097)	0.138 (0.033)	0.132 (0.032)	1.078 (0.058)	0.395 (0.061)	0.380 (0.060)
40%	0.819 (0.035)	0.301 (0.030)	0.278 (0.026)	0.683 (0.082)	0.161 (0.038)	0.154 (0.036)	1.077 (0.041)	0.448 (0.065)	0.428 (0.063)
50%	0.817 (0.029)	0.346 (0.031)	0.325 (0.028)	0.675 (0.084)	0.198 (0.043)	0.188 (0.041)	1.074 (0.036)	0.522 (0.063)	0.495 (0.060)

mean square error (RMSE). Comparison results are listed in Tables 2 and 3. As seen in the tables, we found that both EM and DA substantially outperform MI for all cases. Furthermore, DA imputation exhibits considerable promising accuracy in the prediction of missing values when compared to the EM imputation, especially as the size of observed values becomes small (i.e., missing rate increases).

As another illustration, we attempt to explore classification accuracies between the ML classifier (12) and the Bayesian classifier (21) via PSV. Experimental results in Table 4 indicate that both classifiers are comparable at low-level missing, but Bayesian classifier yields lower misclassification rates as the missing rate increases, though improvements are not substantial.

Finally, we are interested in comparing behaviors of density estimation from ML-fitted and Bayesian posterior predictive aspects. To illustrate this, we use the *salmon*

Table 3: A comparison of prediction accuracies for MI, EM and DA imputations with the standard deviations in parentheses for the *crabs* data set. (Replications=500)

$\eta$	MAE			MARE			RMSE		
	MI	EM	DA	MI	EM	DA	MI	EM	DA
10%	4.063 (0.337)	0.421 (0.055)	0.415 (0.050)	0.202 (0.018)	0.024 (0.003)	0.023 (0.003)	5.391 (0.427)	0.611 (0.114)	0.598 (0.105)
20%	4.008 (0.227)	0.484 (0.041)	0.474 (0.037)	0.200 (0.012)	0.027 (0.002)	0.026 (0.002)	5.343 (0.305)	0.714 (0.090)	0.693 (0.083)
30%	4.037 (0.169)	0.568 (0.044)	0.550 (0.041)	0.202 (0.009)	0.030 (0.002)	0.029 (0.002)	5.384 (0.225)	0.846 (0.096)	0.812 (0.091)
40%	4.036 (0.138)	0.662 (0.044)	0.632 (0.042)	0.203 (0.007)	0.035 (0.002)	0.033 (0.002)	5.381 (0.188)	0.977 (0.092)	0.932 (0.094)
50%	4.039 (0.108)	0.768 (0.052)	0.728 (0.050)	0.202 (0.006)	0.039 (0.002)	0.037 (0.002)	5.386 (0.142)	1.120 (0.102)	1.058 (0.100)

data taken from Johnson and Wichern (2002). This data set has two attributes, the diameter of rings for the first-year freshwater growth and the diameter of rings for the first-year marine growth, for each of 50 Alaskan-born and Canadian born salmon fishes. The ML-fitted density estimation is obtained by plugging the ML estimates into (1). As for Bayesian predictive density, it can be approximated by the use of Rao-Blackwellization

$$\begin{aligned}
p(\mathbf{y}|\mathbf{Y}^o) &= \int p(\mathbf{y}|\mathbf{Y}^o, \Theta)p(\Theta|\mathbf{Y}^o)d\Theta \\
&\approx \frac{1}{L} \sum_{\ell=1}^L p(\mathbf{y}|\Theta^{(\ell)}) \\
&= \frac{1}{L} \sum_{\ell=1}^L \left( \sum_{i=1}^g w_i^{(\ell)} \left( (2\pi)^{-p/2} |\Sigma_i^{(\ell)}|^{-1/2} \exp \left( -\frac{1}{2} \Delta_{ij}^{(\ell)} \right) \right) \right), \quad (22)
\end{aligned}$$

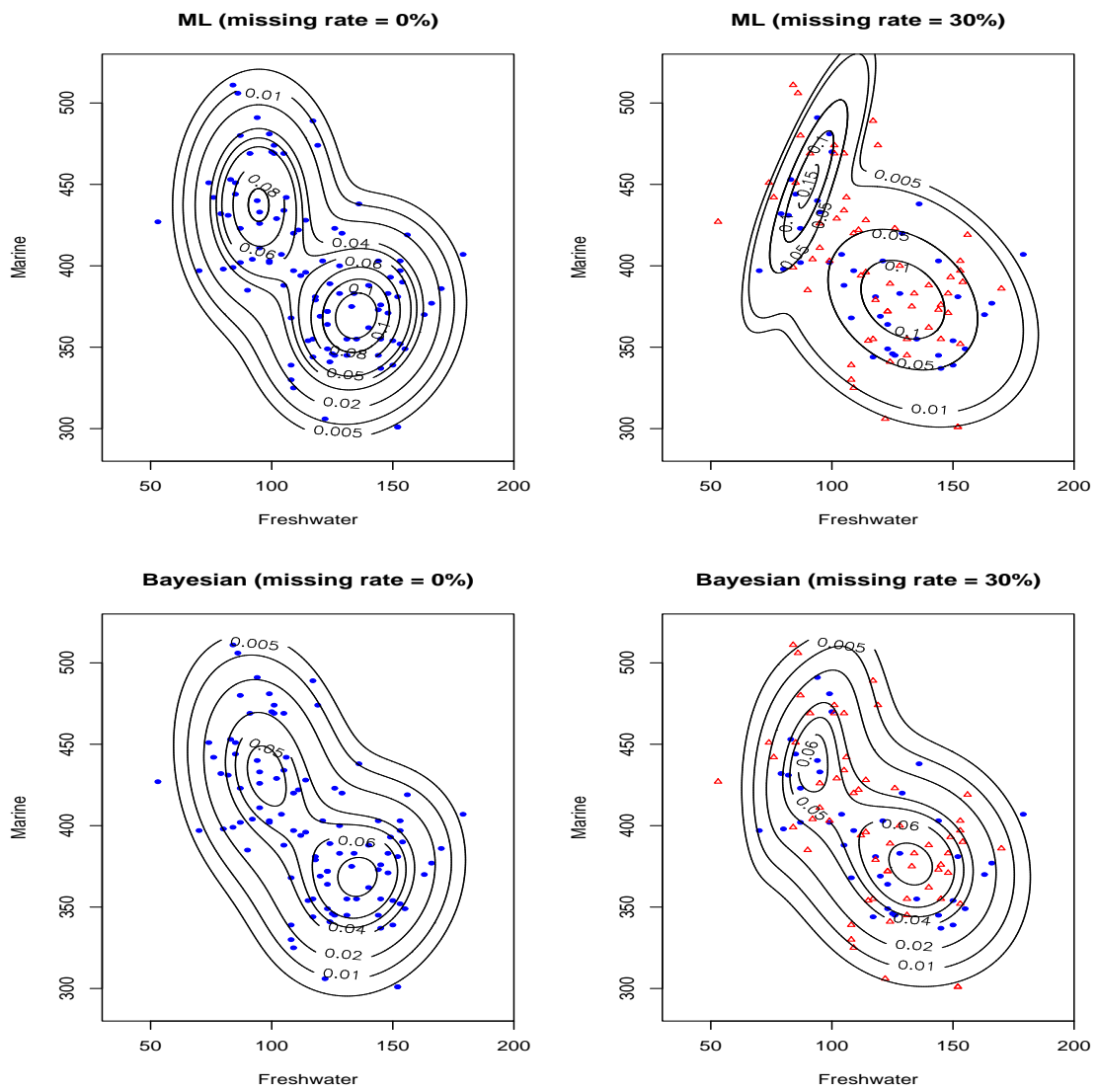
where  $\Delta_{ij}^{(\ell)} = (\mathbf{y} - \boldsymbol{\mu}_i^{(\ell)})^\top \Sigma_i^{(\ell)-1} (\mathbf{y} - \boldsymbol{\mu}_i^{(\ell)})$  and  $\Theta^{(\ell)}$  ( $\ell = 1, \dots, L$ ) is a set of converged Monte Carlo samples generated from the DA algorithm.

Table 4: A comparison of average misclassification rates (%) between ML and Bayesian classifiers. (replicates=500)

$\eta$	<i>Iris</i>		<i>crabs</i>	
	ML	Bayesian	ML	Bayesian
0%	3.33	3.00	7.50	7.30
10%	3.85	3.75	9.75	9.50
20%	5.20	5.00	13.66	13.55
30%	6.90	6.10	19.22	18.80
40%	10.15	9.20	26.75	25.20
50%	13.42	12.30	35.21	33.00

The contour plots obtained by the ML-fitting and Bayesian predictive densities (22) for both completely observed data ( $\eta = 0\%$ ) and both partially observed data ( $\eta = 30\%$ ) are depicted in Figure 1, respectively. Both look similar when data are not missing but using (22) seems to have a relatively smoother appearance. In addition, we found that the ML-fitted contour shapes tend to be distorted at high-level missing and even for moderate-level missing ( $\eta = 30\%$ ). However, the distortion rarely happened while using (22). This indicates that Bayesian learning is more resistant to missing values.

Figure 1: ML and Bayesian density estimation for the two-component salmon data ( $\bullet$ , both attributes are completely observed;  $\triangle$ , one of the two attributes is missing).



## 6. Conclusions

In this paper, two novel EM and DA computational algorithms for learning normal mixture models under a missing information framework are presented. It should be emphasized that our proposed procedures offer neat ways to program with low-cost computation. Experimental results indicate that Bayesian treatment is a worthwhile tool for mixture modelling under a considerable extent of missing information.

Recently, Bayesian and non-Bayesian robust mixture model modelling using the  $t$  distribution has received notable attentions, see Peel and McLanclan (2000), Shoham (2002), Lin et al. (2004) and Wang et al. (2004). Future work will make some kind of comparisons theoretically or empirically among various competitive choices.

## Appendix

### A. Proof of Proposition 2

Suppose  $\mathbf{W} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then for any  $q \times p$  matrix  $\mathbf{A}$  with rank  $q$  ( $q \leq p$ ), we can obtain  $\mathbf{AY} \sim N_p(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$ . With similar arguments, the marginal distributions of  $\mathbf{Y}_j^\circ$  and  $\mathbf{Y}_j^m$  are

$$\begin{aligned} \mathbf{Y}_j^\circ &= \mathbf{O}_j \mathbf{Y}_j \sim \sum_{i=1}^g w_i \phi_{p_j^\circ}(\boldsymbol{\mu}_{ij}^\circ, \boldsymbol{\Sigma}_{ij}^{\circ\circ}), \quad \boldsymbol{\mu}_{ij}^\circ = \mathbf{O}_j \boldsymbol{\mu}_i, \quad \boldsymbol{\Sigma}_{ij}^{\circ\circ} = \mathbf{O}_j \boldsymbol{\Sigma}_i \mathbf{O}_j^\top, \\ \mathbf{Y}_j^m &= \mathbf{M}_j \mathbf{Y}_j \sim \sum_{i=1}^g w_i \phi_{p-p_j^\circ}(\boldsymbol{\mu}_{ij}^m, \boldsymbol{\Sigma}_{ij}^{mm}), \quad \boldsymbol{\mu}_{ij}^m = \mathbf{M}_j \boldsymbol{\mu}_i, \quad \boldsymbol{\Sigma}_{ij}^{mm} = \mathbf{M}_j \boldsymbol{\Sigma}_i \mathbf{M}_j^\top. \end{aligned}$$

Note that the  $\Delta_{ij}$  in (2) can be reexpressed as

$$\Delta_{ij} = \begin{bmatrix} \mathbf{Y}_j^\circ - \boldsymbol{\mu}_{ij}^\circ \\ \mathbf{Y}_j^m - \boldsymbol{\mu}_{ij}^m \end{bmatrix}^\top \begin{bmatrix} \boldsymbol{\Sigma}_{ij}^{\circ\circ} & \boldsymbol{\Sigma}_{ij}^{\circ m} \\ \boldsymbol{\Sigma}_{ij}^{m\circ} & \boldsymbol{\Sigma}_{ij}^{mm} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{Y}_j^\circ - \boldsymbol{\mu}_{ij}^\circ \\ \mathbf{Y}_j^m - \boldsymbol{\mu}_{ij}^m \end{bmatrix}, \quad (23)$$

where  $\Sigma_{ij}^{\text{om}} = \mathbf{O}_j \Sigma_i \mathbf{M}_j^\top$  and  $\Sigma_{ij}^{\text{mo}} = \mathbf{M}_j \Sigma_i \mathbf{O}_j^\top$ . Also, the second and third factors on the right hand side of (23) can be represented by

$$\begin{bmatrix} \Sigma_{ij}^{\text{oo}} & \Sigma_{ij}^{\text{om}} \\ \Sigma_{ij}^{\text{mo}} & \Sigma_{ij}^{\text{mm}} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{I} & -\Sigma_{ij}^{\text{oo}^{-1}} \Sigma_{ij}^{\text{om}} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \Sigma_{ij}^{\text{oo}^{-1}} & \mathbf{0} \\ \mathbf{0} & \Sigma_{ij}^{\text{mm} \cdot \text{o}^{-1}} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\Sigma_{ij}^{\text{mo}} \Sigma_{ij}^{\text{oo}^{-1}} & \mathbf{I} \end{bmatrix},$$

and

$$\begin{bmatrix} \mathbf{Y}_j^{\text{o}} - \boldsymbol{\mu}_{ij}^{\text{o}} \\ \mathbf{Y}_j^{\text{m}} - \boldsymbol{\mu}_{ij}^{\text{m}} \end{bmatrix} = \begin{bmatrix} \mathbf{O}_j (\mathbf{Y}_j - \boldsymbol{\mu}_i) \\ \mathbf{M}_j (\mathbf{Y}_j - \boldsymbol{\mu}_i) \end{bmatrix} = \begin{bmatrix} \mathbf{O}_j \\ \mathbf{M}_j \end{bmatrix} (\mathbf{Y}_j - \boldsymbol{\mu}_i).$$

We then have the following standard results:

$$\begin{aligned} \Sigma_{ij}^{\text{mm} \cdot \text{o}} &= \Sigma_{ij}^{\text{mm}} - \Sigma_{ij}^{\text{mo}} \Sigma_{ij}^{\text{oo}^{-1}} \Sigma_{ij}^{\text{om}} \\ &= \mathbf{M}_j \Sigma_i \mathbf{M}_j^\top - \mathbf{M}_j \Sigma_i \mathbf{O}_j^\top (\mathbf{O}_j \Sigma_i \mathbf{O}_j^\top)^{-1} \mathbf{O}_j \Sigma_i \mathbf{M}_j^\top \\ &= \mathbf{M}_j (\mathbf{I}_p - \Sigma_i \mathbf{S}_{ij}^{\text{oo}}) \Sigma_i \mathbf{M}_j^\top = \mathbf{E}_{ij} \Sigma_i \mathbf{M}_j^\top, \end{aligned}$$

where  $\mathbf{E}_{ij} = \mathbf{M}_j (\mathbf{I}_p - \Sigma_i \mathbf{S}_{ij}^{\text{oo}})$ ,  $\mathbf{S}_{ij}^{\text{oo}} = \mathbf{O}_j^\top (\mathbf{O}_j \Sigma_i \mathbf{O}_j^\top)^{-1} \mathbf{O}_j$ .

Since

$$\begin{aligned} -\Sigma_{ij}^{\text{mo}} \Sigma_{ij}^{\text{oo}^{-1}} \mathbf{O}_j + \mathbf{M}_j &= \mathbf{M}_j - \mathbf{M}_j \Sigma_i \mathbf{O}_j^\top (\mathbf{O}_j \Sigma_i \mathbf{O}_j^\top)^{-1} \mathbf{O}_j \\ &= \mathbf{M}_j (\mathbf{I}_p - \Sigma_i \mathbf{S}_{ij}^{\text{oo}}) \\ &= \mathbf{E}_{ij} \end{aligned}$$

and

$$\begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\Sigma_{ij}^{\text{mo}} \Sigma_{ij}^{\text{oo}^{-1}} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{O}_j \\ \mathbf{M}_j \end{bmatrix} = \begin{bmatrix} \mathbf{O}_j \\ \mathbf{E}_{ij} \end{bmatrix},$$

it suffices to show that

$$\begin{aligned} \boldsymbol{\mu}_{ij}^{\text{m} \cdot \text{o}} &= \boldsymbol{\mu}_{ij}^{\text{m}} + \Sigma_{ij}^{\text{mo}} \Sigma_{ij}^{\text{oo}^{-1}} (\mathbf{Y}_j^{\text{o}} - \boldsymbol{\mu}_{ij}^{\text{o}}) \\ &= \mathbf{M}_j \boldsymbol{\mu}_i + \mathbf{M}_j \Sigma_i \mathbf{O}_j^\top (\mathbf{O}_j \Sigma_i \mathbf{O}_j^\top)^{-1} \mathbf{O}_j (\mathbf{Y}_j - \boldsymbol{\mu}_i) \\ &= \mathbf{M}_j (\boldsymbol{\mu}_i + \Sigma_i \mathbf{S}_{ij}^{\text{oo}} (\mathbf{Y}_j - \boldsymbol{\mu}_i)). \end{aligned}$$

Hence,

$$\begin{aligned} \Delta_{ij} &= (\mathbf{Y}_j^{\text{o}} - \boldsymbol{\mu}_{ij}^{\text{o}})^\top \Sigma_{ij}^{\text{oo}^{-1}} (\mathbf{Y}_j^{\text{o}} - \boldsymbol{\mu}_{ij}^{\text{o}}) + (\mathbf{Y}_j^{\text{m}} - \boldsymbol{\mu}_{ij}^{\text{m} \cdot \text{o}})^\top \Sigma_{ij}^{\text{mm} \cdot \text{o}^{-1}} (\mathbf{Y}_j^{\text{m}} - \boldsymbol{\mu}_{ij}^{\text{m} \cdot \text{o}}) \\ &= (\mathbf{Y}_j - \boldsymbol{\mu}_i)^\top (\mathbf{S}_{ij}^{\text{oo}} + \mathbf{S}_{ij}^{\text{mm} \cdot \text{o}}) (\mathbf{Y}_j - \boldsymbol{\mu}_i) \\ &= \Delta_{ij}^{\text{o}} + \Delta_{ij}^{\text{m} \cdot \text{o}}, \end{aligned}$$

where

$$\begin{aligned}\Delta_{ij}^o &= (\mathbf{Y}_j^o - \boldsymbol{\mu}_{ij}^o)^\top \boldsymbol{\Sigma}_{ij}^{oo^{-1}} (\mathbf{Y}_j^o - \boldsymbol{\mu}_{ij}^o) = (\mathbf{Y}_j - \boldsymbol{\mu}_i)^\top \mathbf{S}_{ij}^{oo} (\mathbf{Y}_j - \boldsymbol{\mu}_i), \\ \Delta_{ij}^{m\cdot o} &= (\mathbf{Y}_j^m - \boldsymbol{\mu}_{ij}^{m\cdot o})^\top \boldsymbol{\Sigma}_{ij}^{mm\cdot o^{-1}} (\mathbf{Y}_j^m - \boldsymbol{\mu}_{ij}^{m\cdot o}) = (\mathbf{Y}_j - \boldsymbol{\mu}_i)^\top \mathbf{S}_{ij}^{mm\cdot o} (\mathbf{Y}_j - \boldsymbol{\mu}_i), \\ \mathbf{S}_{ij}^{mm\cdot o} &= \mathbf{E}_{ij}^\top (\mathbf{E}_{ij} \boldsymbol{\Sigma}_i \mathbf{M}_j^\top)^{-1} \mathbf{E}_{ij}.\end{aligned}$$

Using the fact that  $|\boldsymbol{\Sigma}_i| = |\boldsymbol{\Sigma}_{ij}^{oo}| |\boldsymbol{\Sigma}_{ij}^{mm\cdot o}|$  and above results, we have

$$\begin{aligned}f(\mathbf{Y}_j^m | \mathbf{Y}_j^o) &= \frac{f(\mathbf{Y}_j)}{f(\mathbf{Y}_j^o)} \\ &= \frac{\sum_{i=1}^g w_i \phi_{p_j^o}(\mathbf{Y}_j^o | \boldsymbol{\mu}_{ij}^o, \boldsymbol{\Sigma}_{ij}^{oo}) \phi_{p-p_j^o}(\mathbf{Y}_j^m | \mathbf{Y}_j^o, \boldsymbol{\mu}_{ij}^{m\cdot o}, \boldsymbol{\Sigma}_{ij}^{mm\cdot o})}{\sum_{i=1}^g w_i \phi_{p_j^o}(\mathbf{Y}_j^o | \boldsymbol{\mu}_{ij}^o, \boldsymbol{\Sigma}_{ij}^{oo})} \\ &= \sum_{i=1}^g w_{ij}^* \phi_{p-p_j^o}(\mathbf{Y}_j^m | \mathbf{Y}_j^o, \boldsymbol{\mu}_{ij}^{m\cdot o}, \boldsymbol{\Sigma}_{ij}^{mm\cdot o}),\end{aligned}$$

where  $w_{ij}^* = w_i \phi_{p_j^o}(\mathbf{Y}_j^o | \boldsymbol{\mu}_{ij}^o, \boldsymbol{\Sigma}_{ij}^{oo}) / \sum_{h=1}^g w_h \phi_{p_j^o}(\mathbf{Y}_j^o | \boldsymbol{\mu}_{hj}^o, \boldsymbol{\Sigma}_{hj}^{oo})$ .

## B. Proof of Proposition 3

Letting  $\hat{Z}_{ij}^{(k)} = \mathbb{E}(Z_{ij} | \mathbf{Y}^o, \hat{\boldsymbol{\Theta}}^{(k)})$ ,  $\hat{\boldsymbol{\xi}}_{ij}^{(k)} = \mathbb{E}(Z_{ij} \mathbf{Y}_j | \mathbf{Y}^o, \hat{\boldsymbol{\Theta}}^{(k)})$  and  $\hat{\boldsymbol{\Phi}}_{ij}^{(k)} = \mathbb{E}(Z_{ij} \mathbf{Y}_j \mathbf{Y}_j^\top | \mathbf{Y}^o, \hat{\boldsymbol{\Theta}}^{(k)})$ , we can show that

$$\hat{Z}_{ij}^{(k)} = \Pr(Z_{ij} = 1 | \mathbf{Y}_j^o, \hat{\boldsymbol{\Theta}}^{(k)}) = \frac{\hat{w}_i^{(k)} \phi_{p_j^o}(\mathbf{Y}_j^o | \hat{\boldsymbol{\mu}}_{ij}^{o(k)}, \hat{\boldsymbol{\Sigma}}_{ij}^{oo(k)})}{\sum_{h=1}^g \hat{w}_h^{(k)} \phi_{p_j^o}(\mathbf{Y}_j^o | \hat{\boldsymbol{\mu}}_{hj}^{o(k)}, \hat{\boldsymbol{\Sigma}}_{hj}^{oo(k)})},$$

$$\begin{aligned}\hat{\boldsymbol{\xi}}_{ij}^{(k)} &= \Pr(Z_{ij} = 1 | \mathbf{Y}^o, \hat{\boldsymbol{\Theta}}^{(k)}) \mathbb{E}[\mathbf{Y}_j | Z_{ij} = 1, \mathbf{Y}^o, \hat{\boldsymbol{\Theta}}^{(k)}] \\ &= \mathbb{E}(Z_{ij} | \mathbf{Y}^o, \hat{\boldsymbol{\Theta}}^{(k)}) \mathbb{E}(\mathbf{Y}_j | Z_{ij} = 1, \mathbf{Y}^o, \hat{\boldsymbol{\Theta}}^{(k)}) \\ &= \hat{Z}_{ij}^{(k)} \hat{\mathbf{Y}}_{ij}^{(k)},\end{aligned}$$

and

$$\begin{aligned}\hat{\boldsymbol{\Phi}}_{ij}^{(k)} &= \mathbb{E}(Z_{ij} \mathbf{Y}_j \mathbf{Y}_j^\top | \mathbf{Y}^o, \hat{\boldsymbol{\Theta}}^{(k)}) \\ &= \mathbb{E}(Z_{ij} | \mathbf{Y}^o, \hat{\boldsymbol{\Theta}}^{(k)}) \mathbb{E}(\mathbf{Y}_j \mathbf{Y}_j^\top | Z_{ij} = 1, \mathbf{Y}^o, \hat{\boldsymbol{\Theta}}^{(k)}) \\ &= \hat{Z}_{ij}^{(k)} \left( \hat{\mathbf{Y}}_{ij}^{(k)} \hat{\mathbf{Y}}_{ij}^{(k)\top} + \text{Cov}(\mathbf{Y}_j | Z_{ij} = 1, \mathbf{Y}^o, \hat{\boldsymbol{\Theta}}^{(k)}) \right) \\ &= \hat{Z}_{ij}^{(k)} \left( \hat{\mathbf{Y}}_{ij}^{(k)} \hat{\mathbf{Y}}_{ij}^{(k)\top} + (\mathbf{I}_p - \hat{\boldsymbol{\Sigma}}_i^{(k)} \hat{\mathbf{S}}_{ij}^{oo(k)}) \hat{\boldsymbol{\Sigma}}_i^{(k)} \right).\end{aligned}$$

Since  $\mathbf{Y}_j = \mathbf{O}_j^\top \mathbf{Y}_j^o + \mathbf{M}_j^\top \mathbf{Y}_j^m$  and  $\mathbf{O}_j^\top \mathbf{O}_j (\mathbf{I}_p - \hat{\Sigma}_i \hat{\mathbf{S}}_{ij}^{oo}) = \mathbf{0}$ , we have

$$\begin{aligned}
\hat{\mathbf{Y}}_{ij}^{(k)} &= \mathbb{E}(\mathbf{Y}_j | Z_{ij} = 1, \mathbf{Y}^o, \hat{\Theta}^{(k)}) \\
&= \mathbb{E}(\mathbf{O}_j^\top \mathbf{Y}_j^o + \mathbf{M}_j^\top \mathbf{Y}_j^m | Z_{ij} = 1, \mathbf{Y}^o, \hat{\Theta}^{(k)}) \\
&= \mathbf{O}_j^\top \mathbf{Y}_j^o + \mathbf{M}_j^\top \mathbb{E}(\mathbf{Y}_j^m | Z_{ij} = 1, \mathbf{Y}^o, \hat{\Theta}^{(k)}) \\
&= \mathbf{O}_j^\top \mathbf{Y}_j^o + \mathbf{M}_j^\top \hat{\boldsymbol{\mu}}_{ij}^{m \cdot o(k)} \\
&= \mathbf{O}_j^\top \mathbf{O}_j \mathbf{Y}_j + \mathbf{M}_j^\top \mathbf{M}_j \left( \hat{\boldsymbol{\mu}}_i^{(k)} + \hat{\Sigma}_i^{(k)} \hat{\mathbf{S}}_{ij}^{oo(k)} (\mathbf{Y}_j - \hat{\boldsymbol{\mu}}_i^{(k)}) \right) \\
&= \mathbf{O}_j^\top \mathbf{O}_j \mathbf{Y}_j + (\mathbf{I}_p - \mathbf{O}_j^\top \mathbf{O}_j) \left( \hat{\boldsymbol{\mu}}_i^{(k)} + \hat{\Sigma}_i^{(k)} \hat{\mathbf{S}}_{ij}^{oo(k)} (\mathbf{Y}_j - \hat{\boldsymbol{\mu}}_i^{(k)}) \right) \\
&= \hat{\boldsymbol{\mu}}_i^{(k)} + \hat{\Sigma}_i^{(k)} \hat{\mathbf{S}}_{ij}^{oo(k)} (\mathbf{Y}_j - \hat{\boldsymbol{\mu}}_i^{(k)}) + \mathbf{O}_j^\top \mathbf{O}_j (\mathbf{I}_p - \hat{\Sigma}_i^{(k)} \hat{\mathbf{S}}_{ij}^{oo(k)}) (\mathbf{Y}_j - \hat{\boldsymbol{\mu}}_i^{(k)}) \\
&= \hat{\boldsymbol{\mu}}_i^{(k)} + \hat{\Sigma}_i^{(k)} \hat{\mathbf{S}}_{ij}^{oo(k)} (\mathbf{Y}_j - \hat{\boldsymbol{\mu}}_i^{(k)}),
\end{aligned}$$

and

$$\begin{aligned}
&\text{Cov}(\mathbf{Y}_j | Z_{ij} = 1, \mathbf{Y}^o, \hat{\Theta}^{(k)}) \\
&= \text{Cov}(\mathbf{O}_j^\top \mathbf{Y}_j^o + \mathbf{M}_j^\top \mathbf{Y}_j^m | Z_{ij} = 1, \mathbf{Y}^o, \hat{\Theta}^{(k)}) \\
&= \mathbf{M}_j^\top \text{Cov}(\mathbf{Y}_j^m | Z_{ij} = 1, \mathbf{Y}^o, \hat{\Theta}^{(k)}) \mathbf{M}_j \\
&= \mathbf{M}_j^\top \hat{\Sigma}_{ij}^{mm \cdot o(k)} \mathbf{M}_j \\
&= \mathbf{M}_j^\top \hat{\mathbf{E}}_{ij}^{(k)} \hat{\Sigma}_i^{(k)} \mathbf{M}_j^\top \mathbf{M}_j \\
&= \mathbf{M}_j^\top \mathbf{M}_j (\mathbf{I}_p - \hat{\Sigma}_i^{(k)} \hat{\mathbf{S}}_{ij}^{oo(k)}) \hat{\Sigma}_i^{(k)} \mathbf{M}_j^\top \mathbf{M}_j \\
&= (\mathbf{I}_p - \mathbf{O}_j^\top \mathbf{O}_j) (\mathbf{I}_p - \hat{\Sigma}_i^{(k)} \hat{\mathbf{S}}_{ij}^{oo(k)}) \hat{\Sigma}_i^{(k)} (\mathbf{I}_p - \mathbf{O}_j^\top \mathbf{O}_j) \\
&= (\mathbf{I}_p - \hat{\Sigma}_i^{(k)} \hat{\mathbf{S}}_{ij}^{oo(k)}) \hat{\Sigma}_i^{(k)} (\mathbf{I}_p - \mathbf{O}_j^\top \mathbf{O}_j) \\
&= (\mathbf{I}_p - \hat{\Sigma}_i^{(k)} \hat{\mathbf{S}}_{ij}^{oo(k)}) \hat{\Sigma}_i^{(k)}.
\end{aligned}$$

Therefore,

$$\begin{aligned}
\Omega_{ij}^{(k)} &= \hat{\Phi}_{ij}^{(k)} - 2\hat{\xi}_{ij}^{(k)} \boldsymbol{\mu}_i^\top + \hat{Z}_{ij}^{(k)} \boldsymbol{\mu}_i \boldsymbol{\mu}_i^\top \\
&= \hat{Z}_{ij}^{(k)} \left( \hat{\mathbf{Y}}_{ij}^{(k)} \hat{\mathbf{Y}}_{ij}^{(k)\top} + (\mathbf{I}_p - \hat{\Sigma}_i^{(k)} \hat{\mathbf{S}}_{ij}^{oo(k)}) \hat{\Sigma}_i^{(k)} \right) - 2\hat{Z}_{ij}^{(k)} \hat{\mathbf{Y}}_{ij}^{(k)} \boldsymbol{\mu}_i^\top + \hat{Z}_{ij}^{(k)} \boldsymbol{\mu}_i \boldsymbol{\mu}_i^\top \\
&= \hat{Z}_{ij}^{(k)} \left( (\hat{\mathbf{Y}}_{ij}^{(k)} - \boldsymbol{\mu}_i) (\hat{\mathbf{Y}}_{ij}^{(k)} - \boldsymbol{\mu}_i)^\top + (\mathbf{I}_p - \hat{\Sigma}_i^{(k)} \hat{\mathbf{S}}_{ij}^{oo(k)}) \hat{\Sigma}_i^{(k)} \right).
\end{aligned}$$



## References

- Anderson, E. 1935. The irises of the Gaspé Peninsula, *Bulletin of the American Iris Society*. 59, 2-5.
- Basford K.E., McLachlan G.J. 1985. Estimation of allocation rates in a cluster analysis text. *J. Amer. Statist. Assoc.* 80, 286-293.
- Brooks S.P., Gelman A. 1998. General methods for monitoring convergence of iterative simulations, *J. Comp. Graph. Statist.* 7, 434-455.
- Campbell, N.A., Mahon, R.J. 1974. A multivariate study of variation in two species of rock crab of genus *Leptograpsus*, *Aust. J. Zoology* 22, 417-425.
- Celeux G., Hurn M., Robert C.P. 2000. Computational and inferential difficulties with mixture posterior distributions, *J. Amer. Statist. Assoc.* 95, 957-970.
- Dempster A.P., Laird N.M., Rubin D.B. 1977. Maximum likelihood from incomplete data via the EM algorithm (with discussion), *J. R. Stat. Soc. B.* 39, 1-38.
- Diebolt J., Robert C.P. 1994. Estimation of finite mixture distributions through Bayesian sampling, *J. R. Stat. Soc. B.* 56, 363-375.
- Edwards W.H. Lindman, Savage L.J. 1963. Bayesian statistical inference for psychological research, *Psychol. Rev.* 70, 193-242.
- Escobar M.D., West M. 1995. Bayesian density estimation and inference using mixtures, *J. Amer. Statist. Assoc.* 90, 577-88.
- Fisher R.A. 1936. The use of multiple measurements in taxonomic problems. *Annals of Eugenics.* 7, Part II, 179-188.
- Fruhwirth-Schnatter S. 2001. Markov Chain Monte Carlo estimation of classical and dynamic switching and mixture models, *J. Amer. Statist. Assoc.* 96, 194-209.
- Gelfand A.E., Smith A.F.M., 1990. Sampling based approaches to calculate marginal densities, *J. Amer. Statist. Assoc.* 85, 398-409.

- Geman S., Geman D., 1984. Stochastic relaxation, Gibbs distributions and the Bayesian restoration of images, *IEEE Trans. Pattern Anal. Mach. Intell.* 6, 721-741.
- Ghahramani Z., Jordan M.I., 1994. Supervised learning from incomplete data via an EM approach, In: Cowan, J.D., Tesar, G., Alspector, J. (Eds), *Advances in Neural Information Processing Systems*, vol. 6. Morgan Kaufmann Publishers, San Francisco, CA, pp. 120-127.
- Green P.J. 1995. Reversible jump Markov chain Monte Carlo computation and Bayesian model determination, *Biometrika.* 82, 711-732
- Johnson R.A., Wichern D.W. 2002. *Applied Multivariate Statistical Analysis*, 5th ed, Prentice Hall.
- Lin T.I., Lee J.C., Ni H.F. 2004. Bayesian Analysis of Mixture Modelling using the Multivariate  $t$  Distribution, *Statist. Comput.* 14, 119-130
- Little, R.J.A., Rubin, D.B. 2002. *Statistical analysis with missing data*, 2nd ed. New York, Wiley.
- Liu C.H. 1999. Efficient ML estimation of multivariate normal distribution from incomplete data. *J. Multivariate. Anal.* 69, 206-217.
- McLachlan G.J., Basford K.E. 1988. *Mixture Models: Inference and Application to Clustering*, New York, Marcel Dekker.
- McLachlan G.J., Peel D. 2000. *Finite Mixture Model*, New York, Wiley.
- Peel D., McLachlan G.J. 2000. Robust mixture modeling using the  $t$  distribution, *Statist. Comput.* 10, 339-348.
- Reaven G.M., Miller R.G. 1979. An attempt to define the nature of chemical diabetes using a multidimensional analysis, *Diabetologia.* 16, 17-24.
- Rubin, D.B. 1976. Inference and missing data, *Biometrika.* 63, 581-592.
- Richardson S., Green P.J. 1997. On Bayesian analysis of mixtures with an unknown number of components, *J. R. Stat. Soc. B.* 59, 731-792.

- Schafer J.L. 1997. Analysis of Incomplete Multivariate Data, London, Chapman and Hall.
- Shoham S. 2002. Robust clustering by deterministic agglomeration EM of mixtures of multivariate  $t$ -distributions, Pattern Recognition. 35, 1127-1142.
- Stephens M. 2000. Bayesian analysis of mixture models with an unknown number of components – an alternative to reversible jump methods, Ann. Statist. 28, 40-74.
- Stone M. 1974. Cross-validated choice and assessment of statistical prediction (with discussion). J. R. Stat. Soc. B. 36, 111-147.
- Tanner M.A., Wong W. H. 1987. The calculation of posterior distributions by data augmentation (with discussion), J. Am. Statist. Assoc. 82, 528-550.
- Titterton, D.M., Smith, A.F.M., Markov, U.E. 1985. Statistical Analysis of Finite Mixture Distributions, New York, Wiley.
- Wang H.X., Zhang Q.B., Luo B., Wei S. 2004. Robust mixture modelling using multivariate  $t$  distribution with missing information, Pattern Recognition Lett. 25, 701-710.
- Zhang Z.H., Chan K.L., Wu Y.M., Chen C.B. 2004. Learning a multivariate gaussian mixture model with the reversible jump MCMC algorithm, Statist. Comput. 14, 343-355.