

# Mixtures of Factor Analyzers with Common Factor Loadings for the Clustering and Visualisation of High-Dimensional Data

Jangsun Baek<sup>1</sup> and Geoffrey J. McLachlan<sup>2</sup>

<sup>1</sup>Department of Statistics, Chonnam National University, Gwangju 500-757, South Korea

and

<sup>2</sup>Department of Mathematics and Institute for Molecular Bioscience, University of Queensland, Brisbane QLD 4072, Australia

March, 2008

**Abstract** – Mixtures of factor analyzers enable model-based density estimation and clustering to be undertaken for high-dimensional data, where the number of observations  $n$  is very large relative to their dimension  $p$ . In practice, there is often the need to reduce further the number of parameters in the specification of the component-covariance matrices. To this end, we propose the use of common component-factor loadings, which considerably reduces further the number of parameters. Moreover, it allows the data to be displayed in low-dimensional plots.

**Keywords** – Normal mixture models, mixtures of factor analyzers, common factor loadings, model-based clustering .

# 1 Introduction

Finite mixture models are being increasingly used to model the distributions of a wide variety of random phenomena and to cluster data sets; see, for example, [1]. Let

$$\mathbf{Y} = (Y_1, \dots, Y_p)^T \quad (1)$$

be a  $p$ -dimensional vector of feature variables. For continuous features  $Y_j$ , the density of  $\mathbf{Y}$  can be modelled by a mixture of a sufficiently large enough number  $g$  of multivariate normal component distributions,

$$f(\mathbf{y}; \Psi) = \sum_{i=1}^g \pi_i \phi(\mathbf{y}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i), \quad (2)$$

where  $\phi(\mathbf{y}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$  denotes the  $p$ -variate normal density function with mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ . Here the vector  $\Psi$  of unknown parameters consists of the mixing proportions  $\pi_i$ , the elements of the component means  $\boldsymbol{\mu}_i$ , and the distinct elements of the component-covariance matrices  $\boldsymbol{\Sigma}_i$  ( $i = 1, \dots, g$ ).

The parameter vector  $\Psi$  can be estimated by maximum likelihood. For an observed random sample,  $\mathbf{y}_1, \dots, \mathbf{y}_n$ , the log likelihood function for  $\Psi$  is given by

$$\log L(\Psi) = \sum_{j=1}^n \log f(\mathbf{y}_j; \Psi). \quad (3)$$

The maximum likelihood estimate (MLE) of  $\Psi$ ,  $\hat{\Psi}$ , is given by an appropriate root of the likelihood equation,

$$\partial \log L(\Psi) / \partial \Psi = \mathbf{0}. \quad (4)$$

Solutions of (4) corresponding to local maximizers of  $\log L(\Psi)$  can be obtained via the expectation-maximization (EM) algorithm [2]; see also [3].

Besides providing an estimate of the density function of  $\mathbf{Y}$ , the normal mixture model (2) provides a probabilistic clustering of the observed data  $\mathbf{y}_1, \dots, \mathbf{y}_n$  into  $g$  clusters in terms of their estimated posterior probabilities of component membership of the mixture. The posterior probability  $\tau_i(\mathbf{y}_j; \Psi)$  that the  $j$ th feature vector with observed value  $\mathbf{y}_j$  belongs to the  $i$ th component of the mixture can be expressed by Bayes' theorem as

$$\tau_i(\mathbf{y}_j; \Psi) = \frac{\pi_i \phi(\mathbf{y}_j; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)}{\sum_{h=1}^g \pi_h \phi(\mathbf{y}_j; \boldsymbol{\mu}_h, \boldsymbol{\Sigma}_h)} \quad (i = 1, \dots, g; j = 1, \dots, n). \quad (5)$$

An outright assignment of the data is obtained by assigning each data point  $\mathbf{y}_j$  to the component to which it has the highest estimated posterior probability of belonging.

The  $g$ -component normal mixture model (2) with unrestricted component-covariance matrices is a highly parameterized model with  $d = \frac{1}{2}p(p+1)$  parameters for each component-covariance matrix  $\boldsymbol{\Sigma}_i$  ( $i = 1, \dots, g$ ). Banfield and Raftery [4] introduced a parameterization of the component-covariance matrix  $\boldsymbol{\Sigma}_i$  based on a variant of the standard spectral decomposition of  $\boldsymbol{\Sigma}_i$  ( $i = 1, \dots, g$ ). But if  $p$  is large relative to the

sample size  $n$ , it may not be possible to use this decomposition to infer an appropriate model for the component-covariance matrices. Even if it is possible, the results may not be reliable due to potential problems with near-singular estimates of the component-covariance matrices when  $p$  is large relative to  $n$ .

In this paper, we focus on the use of mixtures of factor analyzers to reduce the number of parameters in the specification of the component-covariance matrices, as discussed in [1, 5, 6]; see also [7]. With the factor-analytic representation of the component-covariance matrices, we have that

$$\boldsymbol{\Sigma}_i = \mathbf{B}_i \mathbf{B}_i^T + \mathbf{D}_i \quad (i = 1, \dots, g), \quad (6)$$

where  $\mathbf{B}_i$  is a  $p \times q$  matrix and  $\mathbf{D}_i$  is a diagonal matrix. As  $\frac{1}{2}q(q-1)$  constraints are needed for  $\mathbf{B}_i$  to be uniquely defined, the number of free parameters in (6) is

$$pq + p - \frac{1}{2}q(q-1). \quad (7)$$

Thus with this representation (6), the reduction in the number of parameters for  $\boldsymbol{\Sigma}_i$  is

$$\begin{aligned} r &= \frac{1}{2}p(p+1) - pq - p + \frac{1}{2}q(q-1) \\ &= \frac{1}{2}\{(p-q)^2 - (p+q)\}, \end{aligned} \quad (8)$$

assuming that  $q$  is chosen sufficiently smaller than  $p$  so that this difference is positive. The total number of parameters is

$$d_1 = (g-1) + 2gp + g\{pq - \frac{1}{2}q(q-1)\}. \quad (9)$$

We shall refer to this approach as MFA (mixtures of factor analyzers).

Even with this MFA approach, the number of parameters still might not be manageable, particularly if the number of dimensions  $p$  is large and/or the number of components (clusters)  $g$  is not small.

In this paper, we therefore consider how this factor-analytic approach can be modified to provide a greater reduction in the number of parameters. We extend the model of [8, 9] to propose the normal mixture model (2) with the restrictions

$$\boldsymbol{\mu}_i = \mathbf{A}\boldsymbol{\xi}_i \quad (i = 1, \dots, g) \quad (10)$$

and

$$\boldsymbol{\Sigma}_i = \mathbf{A}\boldsymbol{\Omega}_i\mathbf{A}^T + \mathbf{D} \quad (i = 1, \dots, g), \quad (11)$$

where  $\mathbf{A}$  is a  $p \times q$  matrix,  $\boldsymbol{\xi}_i$  is a  $q$ -dimensional vector,  $\boldsymbol{\Omega}_i$  is a  $q \times q$  positive definite symmetric matrix, and  $\mathbf{D}$  is a diagonal  $p \times p$  matrix. As to be made more precise in the next section,  $\mathbf{A}$  is a matrix of loadings on  $q$  unobservable factors and its  $p$  columns are taken to be orthonormal; that is,

$$\mathbf{A}^T \mathbf{A} = \mathbf{I}_q, \quad (12)$$

where  $\mathbf{I}_q$  is the  $q \times q$  identity matrix. With the restrictions (10) and (11) on the component mean  $\boldsymbol{\mu}_i$  and covariance matrix  $\boldsymbol{\Sigma}_i$ , respectively, the total number of parameters is reduced to

$$d_2 = (g-1) + p + q(p+g) + \frac{1}{2}(g-1)q(q+1). \quad (13)$$

We shall refer to this approach as MCFA (mixtures of common factor analyzers) since it is formulated via the adoption of a common matrix ( $\mathbf{A}$ ) for the component-factor loadings. We shall show for this approach how the EM algorithm can be implemented to fit this normal mixture model under the constraints (10) and (11). We shall also illustrate how it can be used to provide lower-dimensional plots of the data  $\mathbf{y}_j$  ( $j = 1, \dots, n$ ). It provides an alternative to canonical variates which are calculated from the clusters under the assumption of equal component-covariance matrices.

## 2 Mixtures of Common Factor Analyzers (MCFA)

In this section, we examine the motivation underlying the MCFA approach with its constraints (10) and (11) on the  $g$  component means and covariance matrices  $\boldsymbol{\mu}_i$  and  $\boldsymbol{\Sigma}_i$  ( $i = 1, \dots, g$ ). We shall show that it can be viewed as a special case of the MFA approach.

To see this we first note that the MFA approach with the factor-analytic representation (6) on  $\boldsymbol{\Sigma}_i$  is equivalent to assuming that the distribution of the difference  $\mathbf{Y}_j - \boldsymbol{\mu}_i$  can be modelled as

$$\mathbf{Y}_j - \boldsymbol{\mu}_i = \mathbf{B}_i \mathbf{U}_{ij} + \mathbf{e}_{ij} \quad \text{with prob. } \pi_i \quad (i = 1, \dots, g) \quad (14)$$

for  $j = 1, \dots, n$ , where the (unobservable) factors  $\mathbf{U}_{i1}, \dots, \mathbf{U}_{in}$  are distributed independently  $N(\mathbf{0}, \mathbf{I}_q)$ , independently of the  $\mathbf{e}_{ij}$ , which are distributed independently  $N(\mathbf{0}, \mathbf{D}_i)$ , where  $\mathbf{D}_i$  is a diagonal matrix ( $i = 1, \dots, g$ ).

As noted in the introductory section, this model may not lead to a sufficiently large enough reduction in the number of parameters, particularly if  $g$  is not small. Hence if this is the case, we propose the MCFA approach whereby the distribution of  $\mathbf{Y}_j$  is modelled as

$$\mathbf{Y}_j = \mathbf{A} \mathbf{U}_{ij} + \mathbf{e}_{ij} \quad \text{with prob. } \pi_i \quad (i = 1, \dots, g) \quad (15)$$

for  $j = 1, \dots, n$ , where the (unobservable) factors  $\mathbf{U}_{i1}, \dots, \mathbf{U}_{in}$  are distributed independently  $N(\boldsymbol{\xi}, \boldsymbol{\Omega}_i)$ , independently of the  $\mathbf{e}_{ij}$ , which are distributed independently  $N(\mathbf{0}, \mathbf{D})$ , where  $\mathbf{D}$  is a diagonal matrix ( $i = 1, \dots, g$ ). Here  $\mathbf{A}$  is a  $p \times q$  matrix of factor loadings, which we take to satisfy the relationship (12); that is,  $\mathbf{A}^T \mathbf{A} = \mathbf{I}_q$ .

To see that the MCFA model as specified by (15) is a special case of the MFA approach as specified by (14), we note that we can rewrite (15) as

$$\begin{aligned} \mathbf{Y}_j &= \mathbf{A} \mathbf{U}_{ij} + \mathbf{e}_{ij} \\ &= \mathbf{A} \boldsymbol{\xi}_i + \mathbf{A}(\mathbf{U}_{ij} - \boldsymbol{\xi}_i) + \mathbf{e}_{ij} \\ &= \boldsymbol{\mu}_i + \mathbf{A} \mathbf{K}_i \mathbf{K}_i^{-1} (\mathbf{U}_{ij} - \boldsymbol{\xi}_i) + \mathbf{e}_{ij} \\ &= \boldsymbol{\mu}_i + \mathbf{B}_i \mathbf{U}_{ij}^* + \mathbf{e}_{ij}, \end{aligned} \quad (16)$$

where

$$\boldsymbol{\mu}_i = \mathbf{A} \boldsymbol{\xi}_i, \quad (17)$$

$$\mathbf{B}_i = \mathbf{A} \mathbf{K}_i, \quad (18)$$

$$\mathbf{U}_{ij}^* = \mathbf{K}_i^{-1} (\mathbf{U}_{ij} - \boldsymbol{\xi}_i), \quad (19)$$

and where the  $\mathbf{U}_{ij}^*$  are distributed independently  $N(\mathbf{0}, \mathbf{I}_q)$ . The covariance matrix of  $\mathbf{U}_{ij}^*$  is equal to  $\mathbf{I}_q$ , since  $\mathbf{K}_i$  can be chosen so that

$$\mathbf{K}_i^{-1} \boldsymbol{\Omega}_i \mathbf{K}_i^{-1T} = \mathbf{I}_q \quad (i = 1, \dots, g). \quad (20)$$

On comparing (16) with (14), it can be seen that the MCFA model is a special case of the MFA model with the additional restrictions that

$$\boldsymbol{\mu}_i = \mathbf{A} \boldsymbol{\xi}_i \quad (i = 1, \dots, g), \quad (21)$$

$$\mathbf{B}_i = \mathbf{A} \mathbf{K}_i \quad (i = 1, \dots, g), \quad (22)$$

and

$$\mathbf{D}_i = \mathbf{D} \quad (i = 1, \dots, g). \quad (23)$$

The latter restriction of equal diagonal covariance matrices for the component-specific error terms ( $\mathbf{D}_i = \mathbf{D}$ ) is sometimes imposed with applications of the MFA approach to avoid potential singularities with small clusters (see [5]).

Concerning the restriction (22) that the matrix of factor loadings is equal to  $\mathbf{A} \mathbf{K}_i$  for each component, it can be viewed as adopting common factor loadings before the use of the transformation  $\mathbf{K}_i$  to transform the factors so that they have unit variances and zero covariances. Hence this is why we call this approach mixtures of common factor analyzers. It is also different to the MFA approach in that it considers the factor-analytic representation of the observations  $\mathbf{Y}_j$  directly, rather than the error terms  $\mathbf{Y}_j - \boldsymbol{\mu}_i$ .

As the MFA approach allows a more general representation of the component-covariance matrices and places no restrictions on the component means it is in this sense preferable to the MCFA approach if its application is feasible given the values of  $p$  and  $g$ . If the dimension  $p$  and/or the number of components  $g$  is too large, then the MCFA provides a more feasible approach at the expense of more distributional restrictions on the data. In empirical results some of which are to be reported in the sequel we have found the performance of the MCFA approach is usually at least comparable to the MFA approach for data sets to which the latter is practically feasible. The MCFA approach also has the advantage in that the latent factors in its formulation are allowed to have different means and covariance matrices and are not white noise as with the formulation of the MFA approach. Thus the (estimated) posterior means of the factors corresponding to the observed data can be used to portray the latter in low-dimensional spaces.

### 3 Some Related Approaches

The MCFA approach is similar in form to the approach proposed by Yoshida et al. [8, 9], who also imposed the additional restrictions that the common diagonal covariance matrix  $\mathbf{D}$  of the error terms is spherical,

$$\mathbf{D} = \sigma^2 \mathbf{I}_p, \quad (24)$$

Table 1: The number of parameters for three models

	$p$	$g$	$q$	Number of parameters
MFA	1000	2	2	7999
	1000	4	2	15999
	5000	2	2	39999
	5000	4	2	79999
MCFA	1000	2	2	3008
	1000	4	2	3020
	5000	2	2	15008
	5000	4	2	15020
MCUFSA	1000	2	2	2007
	1000	4	2	2017
	5000	2	2	10007
	5000	4	2	10017

and that the component-covariance matrices of the factors are diagonal. We shall call this approach MUFSA (mixtures of common uncorrelated factor spherical-error analyzers). The total number of parameters with this approach is

$$d_3 = (g - 1) + pq + 1 + 2gq - \frac{1}{2}q(q + 1). \quad (25)$$

In our experience, we have found that these restrictions of sphericity of the errors and of diagonal covariance matrices in the component distributions of the factors can have an adverse effect on the clustering of high-dimensional data sets. The relaxation of these restrictions does considerably increase the complexity of the problem of fitting the model. But we shall show how it can be effected via the EM algorithm with the E- and M-steps being able to be carried out in closed form.

In Table 1, we have listed the number of parameters to be estimated for the MFA, MCFA, and MUFSA models when  $p = 1000, 5000$ ;  $q = 2$ ; and  $g = 2, 4$ . For example, when we cluster  $p = 1000$  dimensional gene expression data into  $g = 2$  groups using  $q = 2$  dimensional factors, the MFA model requires 7999 parameters to be estimated, while the MUFSA needs only 2007 parameters. Moreover, as the number of clusters grows from 2 to 4 the number of parameters for the MFA model grows twice as large as before, but that for MUFSA remains almost the same. As MUFSA needs less parameters to characterize the structure of the clusters, it does not always provide a good fit. It may fail to fit the data unless the directions of the cluster-error distributions are parallel to the axes of the original feature space due the condition of sphericity on the cluster errors. Also, it is assuming that the component-covariance matrices of the factors are diagonal.

Recently, Sanguinetti [10] has considered a method of dimensionality reduction in a cluster analysis context. However, its underlying model assumes sphericity in the specification of the variances/covariances of the factors in each cluster. Our proposed method allows for oblique factors, which provides the extra flexibility needed to cluster more effectively high-dimensional data sets in practice.

## 4 Fitting of Factor-Analytic Models

The fitting of mixtures of factor analyzers as with the MFA approach has been considered in [5], using a variant of the EM algorithm known as the alternating expectation-conditional maximization algorithm (AECM). With the MCFA approach, we have fit to the same mixture model of factor analyzers but with the additional restrictions (10) and (11) on the component means  $\boldsymbol{\mu}_i$  and covariance matrices  $\boldsymbol{\Sigma}_i$ . We also have to impose the restriction (23) of common diagonal covariance matrices  $\mathbf{D}$ . The implementation of the EM algorithm for this model is described in the Appendix. In the EM framework, the component label  $z_j$  associated with the observation  $\mathbf{y}_j$  is introduced as missing data, where  $z_{ij} = (z_j)_i$  is one or zero according as  $\mathbf{y}_j$  belongs or does not belong to the  $i$ th component of the mixture ( $i = 1, \dots, g; j = 1, \dots, n$ ). The unobservable factors  $\mathbf{u}_{ij}$  are also introduced as missing data in the EM framework.

As part of the E-step, we require the conditional expectation of the component labels  $z_{ij}$  ( $i = 1, \dots, g$ ) given the observed data point  $\mathbf{y}_j$  ( $j = 1, \dots, n$ ). It follows that

$$\begin{aligned} E_{\Psi}\{Z_{ij} \mid \mathbf{y}_j\} &= \text{pr}_{\Psi}\{Z_{ij} = 1 \mid \mathbf{y}_j\} \\ &= \tau_i(\mathbf{y}_j; \Psi) \quad (i = 1, \dots, g; j = 1, \dots, n), \end{aligned} \quad (26)$$

where  $\tau_i(\mathbf{y}_j; \Psi)$  is the posterior probability that  $\mathbf{y}_j$  belongs to the  $i$ th component of the mixture. From (5), it can be expressed under the MCFA model as

$$\tau_i(\mathbf{y}_j; \Psi) = \frac{\pi_i \phi(\mathbf{y}_j; \mathbf{A}\boldsymbol{\xi}_i, \mathbf{A}\boldsymbol{\Omega}_i\mathbf{A}^T + \mathbf{D})}{\sum_{h=1}^g \pi_h \phi(\mathbf{y}_j; \mathbf{A}\boldsymbol{\xi}_h, \mathbf{A}\boldsymbol{\Omega}_h\mathbf{A}^T + \mathbf{D})} \quad (27)$$

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

We also require the conditional distribution of the unobservable (latent) factors  $\mathbf{U}_{ij}$  given the observed data  $\mathbf{y}_j$  ( $j = 1, \dots, n$ ). The conditional distribution of  $\mathbf{U}_{ij}$  given  $\mathbf{y}_j$  and its membership of the  $i$ th component of the mixture (that is,  $z_{ij} = 1$ ) is multivariate normal,

$$\mathbf{U}_{ij} \mid \mathbf{y}_j, z_{ij} = 1 \sim N(\boldsymbol{\xi}_{ij}, \boldsymbol{\Omega}_{iy}), \quad (28)$$

where

$$\boldsymbol{\xi}_{ij} = \boldsymbol{\xi}_i + \boldsymbol{\gamma}_i^T(\mathbf{y}_j - \mathbf{A}\boldsymbol{\xi}_i) \quad (29)$$

and

$$\boldsymbol{\Omega}_{iy} = (\mathbf{I}_q - \boldsymbol{\gamma}_i^T \mathbf{A})\boldsymbol{\Omega}_i, \quad (30)$$

and where

$$\boldsymbol{\gamma}_i = (\mathbf{A}\boldsymbol{\Omega}_i\mathbf{A}^T + \mathbf{D})^{-1} \mathbf{A}\boldsymbol{\Omega}_i. \quad (31)$$

We can portray the observed data  $\mathbf{y}_j$  in  $q$ -dimensional space by plotting the corresponding values of the  $\hat{\mathbf{u}}_{ij}$ , which are estimated conditional expectations of the factors  $\mathbf{U}_{ij}$ , corresponding to the observed data points  $\mathbf{y}_j$ . From (28) and (29),

$$\begin{aligned} E(\mathbf{U}_{ij} \mid \mathbf{y}_j, z_{ij} = 1) &= \boldsymbol{\xi}_{ij} \\ &= \boldsymbol{\xi}_i + \boldsymbol{\gamma}_i^T(\mathbf{y}_j - \mathbf{A}\boldsymbol{\xi}_i). \end{aligned} \quad (32)$$

We let  $\hat{\mathbf{u}}_{ij}$  denote the value of the right-hand side of (32) evaluated at the maximum likelihood estimates of  $\boldsymbol{\xi}_i$ ,  $\boldsymbol{\gamma}_i$ , and  $\mathbf{A}$ . We can define the estimated value  $\hat{\mathbf{u}}_j$  of the  $j$ th factor corresponding to  $\mathbf{y}_j$  as

$$\hat{\mathbf{u}}_j = \sum_{i=1}^g \tau_i(\mathbf{y}_j; \hat{\boldsymbol{\Psi}}) \hat{\mathbf{u}}_{ij} \quad (j = 1, \dots, n) \quad (33)$$

where, from (27),  $\tau_i(\mathbf{y}_j; \hat{\boldsymbol{\Psi}})$  is the estimated posterior probability that  $\mathbf{y}_j$  belongs to the  $i$ th component. An alternative estimate of the posterior expectation of the factor corresponding to the  $j$ th observation  $\mathbf{y}_j$  is defined by replacing  $\tau_i(\mathbf{y}_j; \hat{\boldsymbol{\Psi}})$  by  $\hat{z}_{ij}$  in (33), where

$$\hat{z}_{ij} = \arg \max_h \tau_h(\mathbf{y}_j; \hat{\boldsymbol{\Psi}}). \quad (34)$$

## 5 Accuracy of Factor-Analytic Approximations

To illustrate the accuracy of the three factor-analytic approximations as defined above, we performed a small simulation experiment. We generated 100 random vectors from each of  $g = 2$  different three-dimensional multivariate normal distributions. The first distribution had the mean vector  $\boldsymbol{\mu}_1 = (0, 0, 0)^T$  and covariance matrix

$$\boldsymbol{\Sigma}_1 = \begin{pmatrix} 4 & -1.8 & -1 \\ -1.8 & 2 & 0.9 \\ -1 & 0.9 & 2 \end{pmatrix},$$

while the second distribution had mean vector  $\boldsymbol{\mu}_2 = (2, 2, 6)^T$  and covariance matrix

$$\boldsymbol{\Sigma}_2 = \begin{pmatrix} 4 & 1.8 & 0.8 \\ 1.8 & 2 & 0.5 \\ 0.8 & 0.5 & 2 \end{pmatrix}.$$

We applied the MFA, MCFA, and the MCFUSA approaches with  $q = 2$  to cluster the data into two groups. We used the ArrayCluster <http://www.ism.ac.jp/~higuchi/arraycluster.htm>, which was developed by Yoshida *et al.* [9] to implement the MCFUSA approach. There were 2 misclassifications for MFA, 4 for MCFA, and 8 for MCFUSA. As we obtained the parameter estimates for each model we can also predict each observation based on the estimated factor scores and the parameter estimates. In Figures 1, 2, and 3, we have plotted the predicted observations  $\hat{\mathbf{y}}_j$  along with the actual observations  $\mathbf{y}_j$  by the MFA, MCFA, and the MCFUSA approaches. For the MFA approach, the predicted observation is obtained as

$$\hat{\mathbf{y}}_j = \sum_{i=1}^g \tau_i(\mathbf{y}_j; \hat{\boldsymbol{\Psi}}) (\hat{\boldsymbol{\mu}}_i + \hat{\mathbf{B}}_i \hat{\mathbf{u}}_{ij}), \quad (35)$$

where

$$\hat{\mathbf{u}}_{ij} = \hat{\boldsymbol{\alpha}}_i^T (\mathbf{y}_j - \hat{\boldsymbol{\mu}}_i) \quad (36)$$

and

$$\hat{\boldsymbol{\alpha}}_i = (\hat{\mathbf{B}}_i \hat{\mathbf{B}}_i^T + \hat{\mathbf{D}}_i)^{-1} \hat{\mathbf{B}}_i. \quad (37)$$

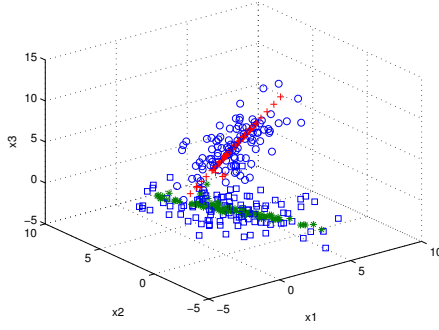


Figure 1: Original observations and the predicted observations by MFA:  $\square$  Group 1;  $\circ$  Group 2; \* predicted for Group 1; + predicted for Group 2

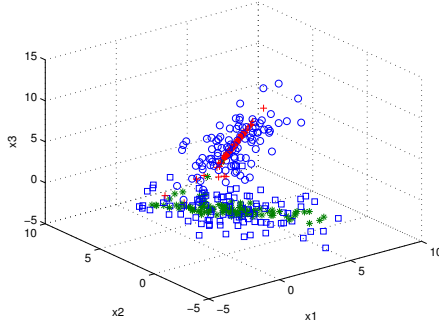


Figure 2: Original observations and the predicted observations by MCFA:  $\square$  Group 1;  $\circ$  Group 2; \* predicted for Group 1; + predicted for Group 2

For the MCFA approach, the predicted observation is

$$\hat{\mathbf{y}}_j = \hat{\mathbf{A}}\hat{\mathbf{u}}_j, \quad (38)$$

where  $\hat{\mathbf{A}}$  is the estimated projection matrix  $\mathbf{A}$  and where  $\hat{\mathbf{u}}_j$  is the estimated factor score for the  $j$ th observation, as defined by (33); similarly, for the MCFSA approach.

The figures show that the original distribution structure of two groups is recovered by the estimated factor scores for MFA and MCFA approaches. Their assumed models are sufficiently flexible to fit the data where the directions of the two cluster-error distributions are not parallel to the axes of the original feature space. On the other hand the predicted observations for the MCFSA approach are not fitted well to the actual distribution of two groups as shown in Figure 3. With this approach, the predicted observations tend to be higher than the actual observations from the first group and lower for those from the second group. This lack of fit is due to the strict assumption of a spherical covariance matrix for each component-error distribution. We measured the difference between the predicted and observed observations by the mean squared error (MSE), where  $\text{MSE} = \sum_{j=1}^{200} (\mathbf{y}_j - \hat{\mathbf{y}}_j)^T (\mathbf{y}_j - \hat{\mathbf{y}}_j) / 200$ . The value of the MSE for the simulated data is 2.30, 3.80, 17.34 for MFA, MCFA, and MCFSA, respectively. As to be expected, the MSE increases in going from MFA to MCFA and then markedly to MCFSA.

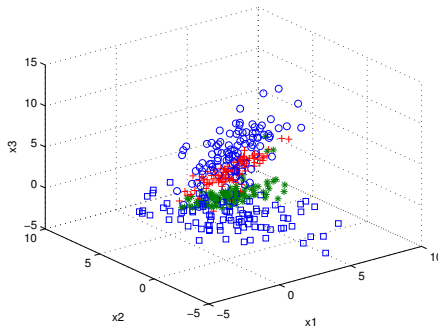


Figure 3: Original observations and the predicted observations by MCFSA:  $\square$  Group 1;  $\circ$  Group 2;  $*$  predicted for Group 1;  $+$  predicted for Group 2

## 6 Applications of MCFA Approach to Clustering

We implemented the MFA, MCFA, and MCFSA approaches to cluster tissues in two different gene expression data sets and individuals in one chemical measurement data set. We compared the consistency between the implied clustering obtained with each approach with the true group membership of each data set. We adopted the clustering corresponding to the local maximizer that gave the largest value of the likelihood as obtained by implementing the EM algorithm for 50 trial starts, comprising 25  $k$ -means starts and 25 random starts. To measure the agreement between a clustering of the data and their true group membership, we used the consistency measures of Jaccard Index [11] and the Adjusted Rand Index (ARI)[12]. Both indices take the value 1 when there is perfect validation between the clustering and the true grouping. The Jaccard Index takes any value between 0 and 1, but the ARI can have negative values.

### 6.1 Example 1: Leukaemia Gene-Expression Data

The first data set concerns the leukaemia tissue samples of Golub et al. [13], in which there are two types of acute leukaemia: Acute Lymphoblastic Leukaemia (ALL) and Acute Myeloid Leukaemia (AML). The data contain the expression levels of 7129 genes on 72 tissues comprising 47 cases of ALL and 25 cases of AML. We preprocessed the data by the following steps: (i) thresholding with a floor of 100 and a ceiling of 16000; (ii) filtering: exclusion of genes with  $\max/\min \leq 5$  and  $(\max - \min) \leq 500$ , where  $\max$  and  $\min$  refer to the maximum and minimum expression levels, respectively, of a particular gene across a tissue sample; (iii) taking logs of the expression levels and then standardizing them across genes to have mean zero and unit standard deviation. Finally, standardizing them across samples to have mean zero and unit standard deviation. This preprocessing resulted in thousand genes being retained. We reduced this set further to 100 genes by selecting genes according to the simple two-sample  $t$ -test as used in [14].

As observed in previous studies, the data on the ALL and AML cases are well separated, confirming the biological distinction between ALL and AML subtypes. Hence all three approaches were able to cluster the into two clusters that almost corresponded

Table 2: Agreement between the clustering result and the true membership of leukaemia data

Similarity Indices	Model	$q$			
		1	2	3	4
Jaccard	MFA	0.951	0.951	0.951	0.951
	MCFA	0.951	0.951	0.951	0.951
	MCUFSA	0.951	0.951	0.951	0.951
ARI	MFA	0.945	0.945	0.945	0.945
	MCFA	0.945	0.945	0.945	0.945
	MCUFSA	0.945	0.945	0.945	0.945

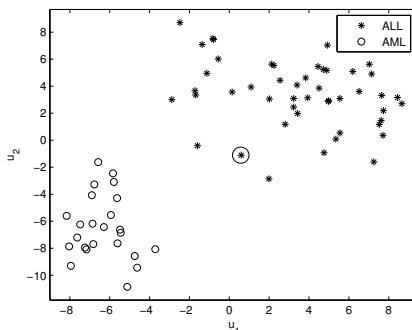


Figure 4: Plot of the factor scores with the MCFA approach

perfectly with the ALL and AML subtypes. There was only one misallocation with the three approaches using either  $q = 1, 2, 3$ , or 4 factors. Hence in Table 2 the values of the Jaccard Index are the same for all three approaches, as are the ARI values.

We plotted the estimated factor scores  $\hat{u}_j = (\hat{u}_{1j}, \hat{u}_{2j})^T$  of the tissues in two-dimensional space according to their clustered membership determined by the MCFA with  $q = 2$  factors (Figure 4). It can be seen from Figure 4 that the two classes ALL and AML are well separated. The one AML tissue that is misallocated as ALL is circled in Figure 4.

## 6.2 Example 2: Paediatric Leukaemia Gene-Expression Data

In the second example, we considered the clustering of the Paediatric Acute Lymphoblastic Leukaemia (ALL) data of Yeoh et al. [15], who analyzed some gene expressions on 7 subtypes of 327 ALL tissues, consisting of BCR-ABL, E2A-PBX1, Hyperdiploid ( $> 50$ ), MLL, Novel, T-ALL, and TEL-AML1 subtypes. They performed an hierarchical cluster analysis of the 327 diagnostic cases using genes selected by several different methods including chi-squared and the  $t$ -statistic. We chose for our analysis the 20 genes with highest chi-squared statistics in each subtype. (This resulted in a total 132 unique genes as some were chosen for more than one subtype. Given that the

Table 3: Agreement between the clustering result and the true membership of paediatric ALL data

Similarity Indices	Model	$q$				
		1	2	3	4	5
Jaccard	MFA	0.632	0.642	0.648	0.555	0.540
	MCFA	NA	0.355	0.597	0.610	0.650
	MCUFSA	-	-	0.593	-	0.584
ARI	MFA	0.720	0.735	0.741	0.650	0.632
	MCFA	NA	0.392	0.685	0.700	0.740
	MCUFSA	-	-	0.687	-	0.677

number of components (subtypes) here is not small with  $g = 7$ , we imposed the constraint (23) of common diagonal matrices  $\mathbf{D}_i$  in the formulation of the MFA approach. This constraint is always imposed with the MCFA approach.

The values of the Jaccard Index and the ARI are listed in Table 3 for the three approaches for each of five levels of the number of factors  $q$  ( $q = 1, 2, 3, 4, 5$ ). In the case of a single factor ( $q = 1$ ), the MFCA approach using  $g = 7$  components clustered the 327 tissues into only 5 clusters, and so the indices were unable to be calculated. This is noted by NA (not available) for  $q = 1$  in Table 3. It can be seen that the performance of the MCFA approach for  $q \geq 3$  factors is comparable to the MFA approach, although their best results for the Jaccard and Adjusted Rand Indices occur for different values of  $q$ . The indices for the latter are greatest for  $q = 3$  factors, whereas they are greatest for the MCFA approach for  $q = 5$  factors. The ArrayCluster program for the implementation of the MUFSA approach gave results for only  $q = 3$  and 5 factors, which are not as high as those for the MCFA approach.

### 6.3 Example 3: Chemical Data with Additional Noise Added

The third example considers the so-called Vietnam data which was considered in Smyth *et al.* [16]. The Vietnam data set consists of the log transformed and standardized concentrations of 17 chemical elements to which four types of synthetic noise variables were added in [16] to study methods for clustering high-dimensional data. We used these data consisting of a total of 67 variables ( $p = 67$ ; 17 chemical concentration variables plus 50 normal noise variables). The concentrations were measured in hair samples from six classes ( $g = 6$ ) of Vietnamese, and the total number of subjects were  $n = 224$ .

The values of the indices for the clustering results for this set are presented in Table 4. It can be seen that the highest value (0.815) of the ARI for the MFCA approach was obtained for  $q = 4$  factors, compared to a highest value of 0.611 for this index at  $q = 3$  factors with the MFA approach. A similar comparison can be made on the basis of the Jaccard Index. Again, it was not possible to obtain results for the MUFSA approach for all values of  $q$ .

Table 4: Agreement between the clustering result and the true membership of Vietnamese data

Similarity Indices	Model	$q$				
		1	2	3	4	5
Jaccard	MFA	0.416	0.507	0.526	0.483	0.425
	MCFA	0.316	0.605	0.590	0.738	0.691
	MCUFSA	-	0.374	0.583	0.577	0.576
ARI	MFA	0.491	0.590	0.611	0.565	0.505
	MCFA	0.331	0.700	0.676	0.815	0.778
	MCUFSA	-	0.440	0.681	0.671	0.675

## 7 Low-Dimensional Plots via MCFA Approach

To illustrate the usefulness of the MFCA approach for portraying the results of a clustering in low-dimensional space, we have plotted in Figure 5 the estimated mean posterior values of the factors  $q_{ij}$  as defined by (33). It can be seen that the clusters are represented in this plot with very little overlap. This is not the case in Figure 6, where the first two canonical variates are plotted. They were calculated using the implied clustering labels. It can be seen from Figure 6 that one cluster is essentially on top of another. The canonical variates are calculated on the basis of the assumption of equal cluster-covariance matrices, which does not apply here. The MCFA approach is not predicated on this assumption and so has more flexibility in representing the data in reduced dimensions.

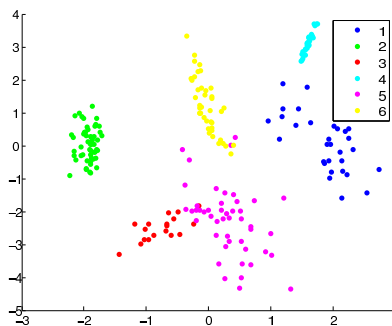


Figure 5: Plot of the (estimated) posterior mean factor scores

## 8 Discussion and Conclusions

In practice, much attention is being given to the use of normal mixture models in density estimation and clustering. However, for high-dimensional data sets, the component-covariance matrices are highly parameterized and some form of reduction in the number of parameters is needed, particularly when the number of observations  $n$  is not large

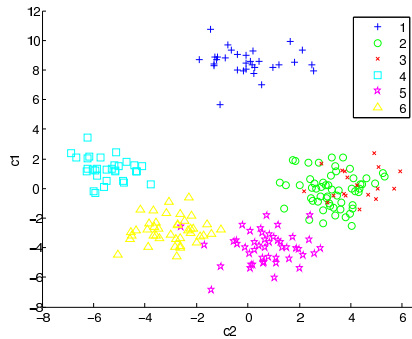


Figure 6: Plot of the first two canonical variates based on the implied clustering via MCFA approach

relative to the number of dimensions  $p$ . One way of proceeding is to work with mixtures of factor analyzers (MFA) as studied in [1, Chapter 8]. This approach achieves a reduction in the number of parameters through its factor-analytic representation of the component-covariance matrices. But it may not provide a sufficient reduction in the number of parameters, particularly when the the number  $g$  of clusters (components) to be imposed on the data is not small. In this paper, we show how in such instances the number of parameters can be reduced appreciably by using a factor-analytic representation of the component-covariance matrices with common factor loadings. The approach is called mixtures of common factor analyzers (MCFA). This sharing of the factor loadings enables the model to be used to cluster high-dimensional into many clusters and to provide low-dimensional plots of the clusters so obtained. The latter plots are given in terms of the (estimated) posterior means of the factors corresponding to the observed data. These projections are not useful with the MFA approach as in its formulation the factors are taken to be white noise with no cluster-specific discriminatory features for the factors.

The MFA approach does allow a more general representation of the component variances/covariances and places no restrictions on the component means. Thus it is more flexible in its modelling of the data. But in this paper we demonstrate that MCFA provides a comparable approach that can be applied in situations where the dimension  $p$  and the number of clusters  $g$  can be quite large. We have presented analyses of both simulated and real data sets to demonstrate the usefulness of the MCFA approach.

In practice, we can use the Bayesian Information Criterion (BIC) of Schwartz [18] to provide a guide to the choice of the number of factors  $q$  and the number of number of components  $g$  to be used. On the latter choice it is well known that regularity conditions do not hold for the usual chi-squared approximation to the asymptotic null distribution of the likelihood ratio test statistic to be valid. However, they do hold for tests on the number of factors at a given level of  $g$ , and so we can also use the likelihood ratio test statistic to choose  $q$ ; see [1, Chapter 8].

## 9 Appendix

The model (15) underlying the MCFA approach can be fitted via the EM algorithm to estimate the vector  $\Psi$  of unknown parameters. It consists of the mixing proportions  $\pi_i$ , the factor component-mean vectors  $\xi_i$ , the distinct elements of the factor component-covariance matrices  $\Omega_i$ , the projection matrix  $\mathbf{A}$  based on sharing of factor loadings, and the common diagonal matrix  $\mathbf{D}$  of the residuals given the factor scores within a component of the mixture. In order to apply the EM algorithm to this problem, we introduce the component-indicator labels  $z_{ij}$ , where  $z_{ij}$  is one or zero according to whether  $\mathbf{y}_j$  belongs or does not belong to the  $i$ th component of the model. We let  $\mathbf{z}_j$  be the component-label vector,  $\mathbf{z}_j = (z_{1j}, \dots, z_{gj})^T$ . The  $\mathbf{z}_j$  are treated as missing data, along with the (unobservable) latent factors  $\mathbf{u}_{ij}$  within this EM framework. The complete-data log likelihood is then given by

$$\log L_c(\Psi) = \sum_{i=1}^g \sum_{j=1}^n z_{ij} \{ \log \pi_i + \log \phi(\mathbf{y}_j; \mathbf{A}\mathbf{u}_{ij}, \mathbf{D}) + \log \phi(\mathbf{u}_{ij}; \xi_i, \Omega_i) \}. \quad (39)$$

### • E-step

On the E-step, we require the conditional expectation of the complete-data log likelihood,  $\log L_c(\Psi)$ , given the observed data  $\mathbf{y} = (\mathbf{y}_1^T, \dots, \mathbf{y}_n^T)^T$ , using the current fit for  $\Psi$ . Let  $\Psi^{(k)}$  be the value of  $\Psi$  after the  $k$ th iteration of the EM algorithm. Then more specifically, on the  $(k+1)$ th iteration the E-step requires the computation of the conditional expectation of  $\log L_c(\Psi)$  given  $\mathbf{y}$ , using  $\Psi^{(k)}$  for  $\Psi$ , which is denoted by  $Q(\Psi; \Psi^{(k)})$ .

We let

$$\tau_{ij}^{(k)} = \tau_i(\mathbf{y}_j; \Psi^{(k)}), \quad (40)$$

where  $\tau_i(\mathbf{y}_j; \Psi)$  is defined by (27). Also, we let  $E_{\Psi^{(k)}}$  refer to the expectation operator, using  $\Psi^{(k)}$  for  $\Psi$ . Then the so-called  $Q$ -function,  $Q(\Psi; \Psi^{(k)})$ , can be written as

$$Q(\Psi; \Psi^{(k)}) = \sum_{i=1}^g \sum_{j=1}^n \tau_{ij}^{(k)} \{ \log \pi_i + w_{1ij}^{(k)} + w_{2ij}^{(k)} \}, \quad (41)$$

where

$$w_{1ij}^{(k)} = E_{\Psi^{(k)}} \{ \log \phi(\mathbf{y}_j; \mathbf{A}\mathbf{u}_{ij}, \mathbf{D}) \mid \mathbf{y}_j, z_{ij} = 1 \} \quad (42)$$

and

$$w_{2ij}^{(k)} = E_{\Psi^{(k)}} \{ \log \phi(\mathbf{u}_{ij}; \xi_i, \Omega_i) \mid \mathbf{y}_j, z_{ij} = 1 \}. \quad (43)$$

### • M-step

On the  $(k+1)$ th iteration of the EM algorithm, the M-step consists of calculating the updated estimates  $\pi_i^{(k+1)}$ ,  $\xi_i^{(k+1)}$ ,  $\Omega_i^{(k+1)}$ ,  $\mathbf{A}^{(k+1)}$ , and  $\mathbf{D}^{(k+1)}$  by solving the equation

$$\partial Q(\Psi; \Psi^{(k)}) / \partial \Psi = \mathbf{0}. \quad (44)$$

The updated estimates of the mixing proportions  $\pi_i$  are given as in the case of the normal mixture model by

$$\pi_i^{(k+1)} = \sum_{j=1}^n \tau_{ij}^{(k)} / n \quad (i = 1, \dots, g). \quad (45)$$

Concerning the other parameters, it can be shown using vector and matrix differentiation that

$$\partial Q(\Psi; \Psi^{(k)}) / \partial \xi_i = \Omega_i^{-1} \sum_{j=1}^n \tau_{ij}^{(k)} E_{\Psi^{(k)}} \{(\mathbf{u}_{ij} - \xi_i) \mid \mathbf{y}_j\}, \quad (46)$$

$$\partial Q(\Psi; \Psi^{(k)}) / \partial \Omega_i^{-1} = \sum_{j=1}^n \tau_{ij}^{(k)} \frac{1}{2} [\Omega_i - E_{\Psi^{(k)}} \{(\mathbf{u}_{ij} - \xi_i)(\mathbf{u}_{ij} - \xi_i)^T \mid \mathbf{y}_j\}], \quad (47)$$

$$\partial Q(\Psi; \Psi^{(k)}) / \partial \mathbf{D}^{-1} = \sum_{i=1}^g \sum_{j=1}^n \tau_{ij}^{(k)} \frac{1}{2} [\mathbf{D} - E_{\Psi^{(k)}} \{(\mathbf{y}_j - \mathbf{A}\mathbf{u}_{ij})(\mathbf{y}_j - \mathbf{u}_{ij})^T \mid \mathbf{y}_j\}], \quad (48)$$

$$\begin{aligned} \partial Q(\Psi; \Psi^{(k)}) / \partial \mathbf{A} = & \sum_{i=1}^g \sum_{j=1}^n \tau_{ij}^{(k)} [\mathbf{D}^{-1} \{\mathbf{y}_j E_{\Psi^{(k)}}(\mathbf{u}_{ij}^T \mid \mathbf{y}_j) \\ & - \mathbf{A} E_{\Psi^{(k)}}(\mathbf{u}_{ij} \mathbf{u}_{ij}^T \mid \mathbf{y}_j)\}]. \end{aligned} \quad (49)$$

On equating (46) to the zero vector, it follows that  $\xi_i^{(k+1)}$  can be expressed as

$$\xi_i^{(k+1)} = \xi_i^{(k)} + \frac{\sum_{j=1}^n \tau_{ij}^{(k)} \gamma_i^{(k)T} (\mathbf{y}_j - \mathbf{A}^{(k)} \xi_i^{(k)})}{\sum_{j=1}^n \tau_{ij}^{(k)}}, \quad (50)$$

where

$$\gamma_i^{(k)} = (\mathbf{A}^{(k)} \Omega_i^{(k)} \mathbf{A}^{(k)T} + \mathbf{D}^{(k)})^{-1} \mathbf{A}^{(k)} \Omega_i^{(k)}. \quad (51)$$

On equating (47) to the null matrix, it follows that

$$\begin{aligned} \Omega_i^{(k+1)} = & \frac{\sum_{j=1}^n \tau_{ij}^{(k)} \gamma_i^{(k)T} (\mathbf{y}_j - \mathbf{A}^{(k)} \xi_i^{(k)}) (\mathbf{y}_j - \mathbf{A}^{(k)} \xi_i^{(k)})^T \gamma_i^{(k)}}{\sum_{j=1}^n \tau_{ij}^{(k)}} \\ & + (\mathbf{I}_q - \gamma_i^{(k)T} \mathbf{A}^{(k)}) \Omega_i^{(k)} \end{aligned} \quad (52)$$

On equating (48) to the zero vector, we obtain

$$\mathbf{D}^{(k+1)} = \text{diag}(\mathbf{D}_1^{(k)} + \mathbf{D}_2^{(k)}), \quad (53)$$

where

$$\mathbf{D}_1^{(k)} = \frac{\sum_{i=1}^g \sum_{j=1}^n \tau_{ij}^{(k)} \mathbf{D}^{(k)} (\mathbf{I}_p - \beta_i^{(k)})}{\sum_{i=1}^g \sum_{j=1}^n \tau_{ij}^{(k)}} \quad (54)$$

and

$$\mathbf{D}_2^{(k)} = \frac{\sum_{i=1}^g \sum_{j=1}^n \tau_{ij}^{(k)} \boldsymbol{\beta}_i^{(k)T} (\mathbf{y}_j - \mathbf{A}^{(k)} \boldsymbol{\xi}_i^{(k)}) (\mathbf{y}_j - \mathbf{A}^{(k)} \boldsymbol{\xi}_i^{(k)})^T \boldsymbol{\beta}_i^{(k)}}{\sum_{i=1}^g \sum_{j=1}^n \tau_{ij}^{(k)}}, \quad (55)$$

and where

$$\boldsymbol{\beta}_i^{(k)} = (\mathbf{A}^{(k)} \boldsymbol{\Omega}_i^{(k)} \mathbf{A}^{(k)T} + \mathbf{D}^{(k)})^{-1} \mathbf{D}^{(k)}. \quad (56)$$

On equating (49) to the null matrix, we obtain

$$\mathbf{A}^{(k+1)} = \left( \sum_{i=1}^g \mathbf{A}_{1i}^{(k)} \right) \left( \sum_{i=1}^g \mathbf{A}_{2i}^{(k)} \right)^{-1}, \quad (57)$$

where

$$\mathbf{A}_{1i}^{(k+1)} = \sum_{j=1}^n \tau_{ij}^{(k)} \{ \mathbf{y}_j \boldsymbol{\xi}_i^{(k)T} + (\mathbf{y}_j - \mathbf{A}^{(k)} \boldsymbol{\xi}_i^{(k)})^T \boldsymbol{\gamma}_i^{(k)} \}, \quad (58)$$

$$\mathbf{A}_{2i}^{(k+1)} = \sum_{j=1}^n \tau_{ij}^{(k)} \{ (\mathbf{I}_q - \boldsymbol{\gamma}_i^{(k)T} \mathbf{A}^{(k)}) \boldsymbol{\Omega}_i^{(k)} + \mathbf{r}_i^{(k)} \mathbf{r}_i^{(k)T} \}, \quad (59)$$

and

$$\mathbf{r}_i^{(k)} = \boldsymbol{\xi}_i^{(k)} + \boldsymbol{\gamma}_i^{(k)T} (\mathbf{y}_j - \mathbf{A}^{(k)} \boldsymbol{\xi}_i^{(k)}) \quad (60)$$

The factor loading matrix  $\mathbf{A}^{(k+1)}$  must be orthogonal in columns, that is,

$$\mathbf{A}^{(k+1)T} \mathbf{A}^{(k+1)} = \mathbf{I}_q. \quad (61)$$

We can use the Cholesky decomposition to find the upper triangular matrix  $\mathbf{C}^{(k+1)}$  of order  $q$  so that

$$\mathbf{A}^{(k+1)T} \mathbf{A}^{(k+1)} = \mathbf{C}^{(k+1)T} \mathbf{C}^{(k+1)}. \quad (62)$$

Then it follows that if we replace  $\mathbf{A}^{(k+1)}$  by

$$\mathbf{A}^{(k+1)} \mathbf{C}^{(k+1)-1}, \quad (63)$$

then it will satisfy the requirement (61). With the adoption of the estimate (63) for  $\mathbf{A}$ , we need to adjust the updated estimates  $\boldsymbol{\xi}_i^{(k+1)}$  and  $\boldsymbol{\Omega}_i^{(k+1)}$  to be

$$\mathbf{C}^{(k+1)} \boldsymbol{\xi}_i^{(k+1)} \quad (64)$$

and

$$\mathbf{C}^{(k+1)} \boldsymbol{\Omega}_i^{(k+1)} \mathbf{C}^{(k)T}, \quad (65)$$

where  $\boldsymbol{\xi}_i^{(k+1)}$  and  $\boldsymbol{\Omega}_i^{(k+1)}$  are given by (50) and (52), respectively.

We have to specify an initial value for the vector  $\boldsymbol{\Psi}$  of unknown parameters in the application of the EM algorithm. A random start is obtained by first randomly assigning the data into  $g$  groups. Let  $n_i$ ,  $\bar{\mathbf{y}}_i$ , and  $\mathbf{S}_i$  be the number of observations, the sample mean, and the sample covariance matrix, respectively, of the  $i$ th group of the data so obtained ( $i = 1, \dots, g$ ). We then proceed as follows:

- Set  $\pi_i^{(0)} = n_i/n$ .
- Generate random numbers from the standard normal distribution  $N(0, 1)$  to obtain values for the  $(j, k)$ th element of  $\mathbf{A}^*$  ( $j = 1, \dots, p; k = 1, \dots, q$ ).
- Implement the Cholesky decomposition so that  $\mathbf{A}^{*T} \mathbf{A}^* = \mathbf{C}^T \mathbf{C}$ , where  $\mathbf{C}$  is the upper triangle matrix of order  $q$ , and define  $\mathbf{A}^{(0)}$  by  $\mathbf{A}^* \mathbf{C}^{-1}$ .
- On noting that the transformed data  $\mathbf{D}^{-1/2} \mathbf{Y}_j$  satisfies the probabilistic PCA model of Tipping and Bishop [17] with  $\sigma_i^2 = 1$ , it follows that for a given  $\mathbf{D}^{(0)}$  and  $\mathbf{A}^{(0)}$ , we can specify  $\Sigma_i^{(0)}$  as

$$\Omega_i^{(0)} = \mathbf{A}^{(0)T} \mathbf{D}^{(0)1/2} \mathbf{H}_i (\Lambda_i - \tilde{\sigma}_i^2 I_q) \mathbf{H}_i^T \mathbf{D}^{(0)1/2} \mathbf{A}^{(0)},$$

where  $\tilde{\sigma}_i^2 = \sum_{h=q+1}^p \lambda_{ih}/(p-q)$ . The  $q$  columns of the matrix  $\mathbf{H}_i$  are the eigenvectors corresponding to the eigenvalues  $\lambda_{i1} \geq \lambda_{i2} \geq \dots \geq \lambda_{iq}$  of

$$\mathbf{D}^{(0)-1/2} \Omega_i^{(0)} \mathbf{D}^{(0)-1/2}, \quad (66)$$

where  $\Omega_i^{(0)}$  is the covariance matrix of the  $\mathbf{y}_j$  in the  $i$ th group, and  $\Lambda_i$  is the diagonal matrix with diagonal elements equal to  $\lambda_{i1}, \dots, \lambda_{iq}$ . Concerning the choice of  $\mathbf{D}^{(0)}$ , we can take  $\mathbf{D}^{(0)}$  to be the diagonal matrix formed from the diagonal elements of the (pooled) within-cluster sample covariance matrix of the  $\mathbf{y}_j$ . The initial value for  $\xi_i$  is  $\xi_i^{(0)} = \mathbf{A}^{(0)T} \bar{\mathbf{y}}_i$ .

Some clustering procedure such as  $k$ -means can be used to provide non-random partitions of the data, which can be used to obtain another set of initial values for the parameters. In our analyses we used both initialization methods.

## Acknowledgement

The work of J. Baek was supported by the Korea Research Foundation Grant funded by the Korean Government (MOEHRD, Basic Research Promotion Fund, KRF-2007-521-C00048). The work of G. McLachlan was supported by the Australian Research Council.

## REFERENCES

- [1] G.J. McLachlan and D. Peel, *Finite Mixture Models*. New York: Wiley, 2000.
- [2] A.P. Dempster, N.M. Laird, and D.B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm (with discussion)," *Journal of the Royal Statistical Society: Series B*, vol. 39, pp. 1–38, 1977.
- [3] G.J. McLachlan and T. Krishnan, *The EM Algorithm and Extensions*. Second Edition. New York: Wiley, 2008.

- [4] J.D. Banfield and A.E. Raftery, “Model-based Gaussian and non-Gaussian clustering,” *Biometrics*, Vol. 49, pp. 803–821, 1993.
- [5] G.J. McLachlan, D. Peel, and R.W. Bean, “Modelling high-dimensional data by mixtures of factor analyzers,” *Computational Statistics & Data Analysis*, vol. 41, pp. 379–388, 2003.
- [6] G.J. McLachlan, R.W. Bean, and L. Ben-Tovim Jones, “Extension of the mixture of factor analyzers model to incorporate the multivariate  $t$  distribution,” *Computational Statistics & Data Analysis*, Vol. 51, 5327–5338, 2007.
- [7] G.E. Hinton, P. Dayan, and M. Revow, M. “Modeling the manifolds of images of handwritten digits,” *IEEE Transactions on Neural Networks*, vol. 8, pp. 65–73, 1997.
- [8] R. Yoshida, T. Higuchi, and S. Imoto, “A mixed factors model for dimension reduction and extraction of a group structure in gene expression data,” *Proceedings of the 2004 IEEE Computational Systems Bioinformatics Conference*, pp. 161–172, 2004.
- [9] R. Yoshida, T. Higuchi, S. Imoto, and S. Miyano, “ArrayCluster: an analytic tool for clustering, data visualization and model finder on gene expression profiles,” *Bioinformatics*, vol. 22, pp. 1538–1539, 2006.
- [10] G. Sanguinetti, “Dimensionality reduction of clustered data sets,” *IEEE Transactions Pattern Analysis and Machine Intelligence*, Vol. 30, pp. 535–540.
- [11] P. Jaccard, “Distribution de la florine alpine dans la Bassin de Dranses et dans quelques regions voisines,” *Bulletin de la Société Vaudoise des Sciences Naturelles*, Vol. 37, pp. 241–272, 1901.
- [12] L. Hubert and P. Arabie, “Comparing partitions,” *Journal of Classification*, Vol. 2, 193–218, 1985.
- [13] T.R. Golub, D.K. Slonim, P. Tamayo, C. Huard, M. Gaasenbeek, P. Mesirov, H. Coller, M.L. Loh, J.R. Downing, M.A. Caligiuri, *et al.*, “Molecular classification of cancer: Class discovery and class prediction by gene expression monitoring,” *Science*, 286, 531-537, 1999.
- [14] D. Nguyen, and D. Roake, “Tumor classification by partial least squares using microarray gene expression data,” *Bioinformatics*, Vol. 18, pp. 39–50, 2002.
- [15] E. Yeoh, and Ross, M.E., *et al.*, “Classification, subtype discovery, and prediction of outcome in pediatric acute lymphoblastic leukemia by gene expression profiling,” *Cancer Cell*, Vol. 1, pp. 133–143, 2002.
- [16] C. Smyth, D. Coomans, and Y. Everingham, “Clustering noisy data in a reduced dimension space via multivariate regression trees,” *Pattern Recognition*, Vol. 39, pp. 424–431, 2006.
- [17] M.E. Tipping, and C.M. Bishop, “Mixtures of probabilistic principal component analysers,” *Neural Computation*, Vol. 11, pp. 443–482, 1999.
- [18] G. Schwarz, “Estimating the dimension of a model,” *Annals of Statistics*, Vol. 6, pp. 461–464, 1978.