

## MIXTURE MODEL CLUSTERING USING THE MULTIMIX PROGRAM

LYNETTE HUNT<sup>1</sup> AND MURRAY JORGENSEN\*<sup>1</sup>

University of Waikato

### Summary

Hunt (1996) implemented the finite mixture model approach to clustering in a program called MULTIMIX. The program is designed to cluster multivariate data that have categorical and continuous variables and that possibly contain missing values. This paper describes the approach taken to design MULTIMIX and how some of the statistical problems were dealt with. As an example, the program is used to cluster a large medical dataset.

**Key words:** cluster analysis; EM algorithm; latent class analysis; local independence; multivariate normal distribution; location model; prostate cancer data.

### 1. Introduction

This paper is concerned with the statistical analysis of multivariate data from a mixture of finitely many populations when there is no information about membership in any component population. This is known as *cluster analysis* or *unsupervised learning*. The goal is to partition the sample into groups so that members of a group are as similar as possible. This is usually done by any one of a number of deterministic algorithms, the most common of which we discuss below.

Cluster analysis is different from *discriminant analysis* where it is possible to classify members of a random sample from a mixture of populations according to which population they come from; see Mardia, Kent & Bibby (1979 Chap. 11) and McLachlan (1992).

There are many different methods for cluster analysis. They can be broadly categorized as hierarchical or non-hierarchical. Clustering, using hierarchical methods, is generally obtained through either agglomerative algorithms, that begin with a cluster for every observation and successively merge clusters, or divisive algorithms, that begin with a single cluster and continually split clusters.

It is possible to visualize two extremes, one in which each object is considered to be a single member cluster, and one in which all  $n$  objects are contained in a single cluster. Each cluster obtained at any stage in the procedure is a combination or division of clusters at other stages. A hierarchical strategy finds an efficient path between these two extremes.

Once an object has been assigned to a cluster under a hierarchical strategy, there is no provision for reallocation of the objects that have been poorly allocated at an earlier stage in the process. Each stage of the analysis involves the computation of the cluster similarity (or distance) matrix. Since the clusters at any stage are obtained by the fusion (agglomerative methods), or division (divisive methods) of clusters from the previous stage, these methods

---

Received March 1997; revised November 1998; accepted November 1998.

\* Author to whom correspondence should be addressed.

<sup>1</sup> Dept of Statistics, University of Waikato, Private Bag 3105, Hamilton, New Zealand.  
email: maj@waikato.ac.nz

**Acknowledgments.** We are grateful to Kaye Basford for making her programs available to us in electronic form.

lead to a hierarchical structure of the objects. This is represented by a dendrogram, also known as a tree diagram.

### 1.1. Similarity Matrix Clustering Techniques

Hierarchical clustering techniques are usually implemented with the data represented by a matrix of proximities ( $d_{ij}$ ), where  $d_{ij}$  is the proximity of observations  $i$  and  $j$ . The proximity  $d_{ij}$  can be either a similarity or a dissimilarity measure. To convert a dissimilarity into a similarity index we may, for example, divide it by the greatest dissimilarity observed in the data and subtract this from 1.

Proximities may be obtained in various ways, one method being to ask a number of people to subjectively assess all pairs of observations in a small set for degree of similarity, recording the answer as a number between 0 (least similar) and 1 (most similar). The similarities for analysis can then be obtained by averaging the subjective similarities over the panel of judges.

More often, each observation has a number of measured attributes or variables, often at differing levels of measurement (binary, nominal, ordinal, interval or ratio), and we require some means of calculating proximities from the data. Anderberg (1973 Chaps 4, 5) and Gordon (1981 Chap. 2) give surveys of many methods of calculating proximities for the case of a single variable. Typical examples are the Euclidean distance for interval variables and the Jacard coefficient  $n_{11}/(n_{11} + n_{10} + n_{01})$  for binary (0,1) data. Another binary coefficient is the simple matching coefficient  $(n_{11} + n_{00})/(n_{11} + n_{10} + n_{01} + n_{00})$ ; indeed, Anderberg lists 14 possibilities, though deprecating five of these. Which notion of proximity makes the most sense depends on subject area considerations.

Once similarity measures  $s_{ijk}$  comparing observation  $i$  with observation  $j$  have been selected for each attribute  $k$  they may be combined, essentially by averaging over the attributes. In the case of a rare binary attribute  $k$  we may wish to exclude  $s_{ijk}$  from the average as being uninformative about the similarity of  $i$  and  $j$ . Gower (1971) gives details about combining similarity measures in this way.

Agglomerative hierarchical techniques differ primarily in how they measure the distance or similarity of two clusters, where a cluster may, at times, consist of a single observation only. For example, the Euclidean distance  $d_{ij}$  between the two observations  $x_i$  and  $x_j$  is defined as  $d_{ij} = [(\mathbf{x}_i - \mathbf{x}_j)'(\mathbf{x}_i - \mathbf{x}_j)]^{1/2}$ , while the Mahalanobis distance is defined as  $d_{ij} = [(\mathbf{x}_i - \mathbf{x}_j)' \hat{\Sigma}^{-1}(\mathbf{x}_i - \mathbf{x}_j)]^{1/2}$ , where  $\hat{\Sigma}$  is the within cluster covariance matrix. Mardia *et al.* (1979) and Gordon (1981) give further details on the properties of these distances and other distance measures.

In single linkage (nearest neighbour) clustering, the distance between two clusters is defined as the distance between their two nearest neighbours:

$$d_{AB} = \min_{i \in A, j \in B} (d_{ij}),$$

where  $d_{AB}$  is the dissimilarity between two clusters  $A$  and  $B$  and  $d_{ij}$  is the dissimilarity between two observations  $i$  and  $j$ . This technique can lead to 'rod' type elongated clusters.

With complete linkage (farthest neighbour) clustering, the distance between two clusters is defined as the distance between their two farthest neighbours:

$$d_{AB} = \min_{i \in A, j \in B} (d_{ij}),$$

This method tends to produce compact clusters.

Other standard linkage methods replace the 'min' and the 'max' of the above methods by measures of central tendency.

Lance & Williams (1967) give a general agglomerative algorithm with which many of the common hierarchical linkage methods can be described. If two clusters  $R$  and  $S$  amalgamate to form a new cluster  $T$ , the dissimilarity between this cluster and any other cluster can be expressed in an equation form. Gordon (1981) includes a table of the algorithm parameters for different techniques.

With hierarchical clustering, the number of clusters is obtained by selecting one of the clusterings in the nested sequence of groupings displayed in the dendrogram. The most common method used is to examine the dendrogram for large changes in the distance or dissimilarity between adjacent fusion levels. A 'large' change when going from  $K$  to  $K - 1$  groups might be indicative of  $K$  groups. This criterion is somewhat subjective. Other sources of subjectivity lie in the choice of similarity metric for each attribute and the choice of linkage method.

These methods are widely implemented in statistical packages and can be useful for preliminary exploration of small multivariate datasets, especially in combination with visualization techniques such as a plot of the first two principal component scores and a Minimal Spanning Tree (Gower & Ross, 1969). They are less satisfactory with large datasets (hundreds rather than tens) because of the large number of pairwise similarities which must be processed, and because of the enhanced possibilities for unfortunate and irreversible amalgamations of clusters at an early stage.

An important problem with the use of these forms of cluster analysis lies in the many ways in which the subjective decisions made by the analyst may influence the outcome. The analyst must choose

- (i) the form of the proximity index,
- (ii) the linkage method, and
- (iii) the similarity level at which to 'cut' the dendrogram, or equivalently, the number of groups.

## 1.2. Optimization Based Techniques

Nonhierarchical techniques of cluster analysis have the same extremes as hierarchical techniques, that is,  $n$  clusters consisting of one observation and one cluster with all  $n$  observations in it. However, nonhierarchical techniques allow points to be reallocated to other clusters during the clustering process. These techniques of cluster analysis often use optimization procedures in which observations are transferred between clusters with the aim of optimizing some clustering criterion that rewards both within-cluster similarity and between-cluster differences. Once again, there are many methods available because of various optimizing criteria and different optimizing algorithms; see Everitt (1980) and Hand (1981) for further discussions on these procedures.

The  $k$ -means algorithm of Hartigan (1975) is a commonly used optimization technique. The means of each of the  $k$  initial clusters are found, and then each data point is examined to see if it is closer to the mean of another cluster than to the mean of its current cluster. If it is, the point is transferred and the cluster means are recalculated. The means can be recalculated after each data point has been reallocated, or after all the data points have been examined and those that needed reallocation have been transferred. The means of the  $k$  clusters are

calculated and the process is repeated. In this procedure, the cluster mean is the point that minimizes the sum of squares of the distances (to that point) of the observations in that cluster.

The 'classification likelihood' approach is a nonhierarchical technique that uses a form of likelihood function as a clustering criterion. This approach applies a probabilistic formulation in which the (vector) observations  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are each assumed to arise from any one of  $K$  possible subpopulations with probability density function  $f(\mathbf{x}; \boldsymbol{\theta}_k)$  for  $k = 1, \dots, K$ , where  $K$  is fixed. This approach differs from the discriminant analysis problem because it is not known which subpopulation the observation comes from. Let

$$z_{ik} = \begin{cases} 1 & \text{if observation } i \in \text{subpopulation } k; \\ 0 & \text{if observation } i \notin \text{subpopulation } k, \end{cases}$$

and define the vector of indicator variables as  $\mathbf{z}_i = (z_{i1}, \dots, z_{iK})'$ . The likelihood function is

$$L_{\text{Class}}(\mathbf{z}_1, \dots, \mathbf{z}_n, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K) = \prod_{i=1}^n \prod_{k=1}^K \{f(\mathbf{x}_i; \boldsymbol{\theta}_k)\}^{z_{ik}}.$$

Let  $\mathbf{z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)$  and  $\boldsymbol{\phi} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K)$ . Maximization of  $L_{\text{Class}}(\mathbf{z}, \boldsymbol{\phi})$ , the likelihood for the complete data is with respect to  $\boldsymbol{\phi}$  and  $\mathbf{z}$ . That is, the unobservable indicator variables  $\mathbf{z}_1, \dots, \mathbf{z}_n$  are treated as unknown parameters to be estimated along with  $\boldsymbol{\phi}$ . The maximization process can be carried out by computing the maximum value of the likelihood over all possible partitions of the  $n$  observations to the  $K$  groups. Several authors including Scott & Symons (1971), Sclove (1977) and Symons (1981) considered this approach. More recently Banfield & Raftery (1993) extended the methods of Scott & Symons (1971) and their approach is discussed below (Section 6.3). Unfortunately, with this procedure, the  $z_{ij}$  increase in number with the number of observations, and the maximum likelihood (ML) estimates are not consistent (McLachlan & Basford, 1988).

Using the classification likelihood approach, Scott & Symons (1971) showed that the assumption that  $\mathbf{x}_i \sim N(\boldsymbol{\mu}_k, \boldsymbol{\Sigma})$  for  $k = 1, \dots, K$ , led to the cluster analysis procedure based on minimizing  $|W|$ , the determinant of the pooled within-group dispersion matrix. Friedman & Rubin (1967) discussed this method of cluster analysis. Scott & Symons (1971) found that this approach has the tendency to divide the data into clusters of equal size if the separation between the subpopulations is not large. Marriot (1975) pointed out that the maximum (classification) likelihood estimates are not consistent under the assumption of underlying normal distributions with a common covariance structure. Bryant & Williamson (1978) showed that the approach can also be expected to give biased results. Binder (1978) and Symons (1981) gave Bayesian versions of this method.

Although usually considered to be nonhierarchical clustering techniques, criterion optimization methods may be used in a hierarchical fashion by applying the algorithm repeatedly to subdivide clusters found earlier. Such an approach usually leads to clusters that are not themselves optimal on the criterion.

### 1.3. Clustering Methods Based on Finite Mixture Models

A vast quantity of literature is available on algorithmic cluster analysis. For comprehensive reviews of clustering techniques see Cormack (1971), Jardine & Sibson (1971), Everitt (1980) and Gordon (1981). For clustering algorithms see Hartigan (1975) and James (1985).

There are some inescapable drawbacks shared by all these traditional approaches to clustering: any randomness in the sample is not reflected and small perturbations in the sample may lead to quite different groups being formed. Further, experience with real mixed populations

shows that quite often there is substantial overlapping, whereas by design most traditional clustering algorithms tend to come up with compact non-overlapping clusters.

An alternative to algorithmic cluster analysis is to adopt a statistical formulation similar to that of discriminant analysis, and regard the observations to be clustered as a random sample from a finite mixture of distributions. However, unlike discriminant analysis, the observations are not identified as belonging to a particular group, and there is often very little information about the form of the population distributions for each group. By making generic distributional assumptions we obtain a well specified model, whose parameters can be estimated by the method of maximum likelihood. The estimated conditional probabilities of group membership can be estimated by Bayes' rule using the parameter estimates. These probabilities can be used when the algorithm has converged, to obtain a probabilistic assignment of observations to clusters. Furthermore, the estimated component distributions together with the estimated proportions for each component provide a concise description of what may be a very complicated set of data.

As with any clustering method, clustering by finite mixture models also imposes a structure on the data. It is possible to check the overall fit of the mixture model to the data, although the individual components cannot be checked unless the clusters turn out to be well separated. The mixture likelihood approach can be seen as an example of a nonhierarchical clustering technique. But a unification with the mainstream of statistical modelling is achieved because clustering methods based on mixture models allow estimation and hypothesis testing within the framework of standard statistical theory (Aitkin, Anderson & Hinde, 1981).

## 2. Earlier Work in Mixture Model Clustering

The MULTIMIX program to be described later in this paper builds on earlier approaches and is most easily understood as an extension and unification of some of these. The estimation problem for finite mixtures of normal distributions has quite a lengthy history.

### 2.1. Mixtures of Normal Distributions

Karl Pearson put forward a solution in the case of a mixture of two univariate distributions with unequal variances, using the method of moments (Pearson, 1894). This was a difficult problem and involved the solution of a ninth degree polynomial equation. Later investigation showed that likelihood estimation was more efficient than the method of moments for this problem (Tan & Chan, 1972).

Maximum likelihood estimation for the parameters in mixture distributions was suggested by Rao (1948), who used Fisher's method of scoring for the estimation of parameters in a mixture of two univariate normal distributions with equal variances. This appeared to be the first use of likelihood estimation for mixtures (Everitt & Hand, 1981). However, Butler (1986) notes that there was an investigation by Newcomb (1886) of the maximum likelihood estimation of the parameters of a mixture of  $K$  univariate normal populations with known variances. His investigation could be interpreted as an application of the EM algorithm (Dempster, Laird & Rubin, 1977). Butler also found that Jeffreys (1932) had essentially used the EM algorithm to compute the estimates of the means in two univariate normal populations, which had known variances and which were mixed in unknown proportions.

With the advent of high speed computers, interest increased in the likelihood estimation of the parameters of mixture distributions. Hasselblad (1966, 1969) applied ML estimation for the parameters of a mixture of  $K$  univariate normal distributions with equal variances,

and then for mixtures of distributions from the exponential family. Day (1969) estimated the components of a mixture of two multivariate normal distributions with equal covariances. Wolfe (1967, 1970) used ML estimation for the parameters of a mixture of  $K$  multivariate normal distributions with unequal covariances, and also a mixture of Bernoulli distributions. These three researchers all presented their solutions in iterative forms that could be viewed as applications of the EM algorithm.

For additional references on finite mixtures, see the monographs on finite mixture distributions by Everitt & Hand (1981), Titterton, Smith & Makov (1985), and McLachlan & Basford (1988), the reviews by Redner & Walker (1984) and the encyclopedia entry by Everitt (1985).

## 2.2. Basford's Mixture-fitting Programs

In their monograph on mixture models and their application to clustering, McLachlan & Basford (1988) focus on the use of  $p$ -variate normal distributions for the component models and consider mainly continuous variables. Their book includes listings of the FORTRAN 66 source code for four programs that estimate the parameters of normal mixture models in various situations. The program of most relevance for cluster analysis is KMM, which fits a mixture of multivariate normal distributions, with either arbitrary or common covariance matrices, by maximum likelihood using the EM algorithm. In designing the MULTIMIX program, we sought to extend and modify KMM to enhance its suitability as a general purpose nonhierarchical clustering program. In fact, MULTIMIX was written from scratch, but its output was tested against that from KMM where possible.

Some development beyond KMM was necessary because of three major difficulties which frustrate the application of multivariate normal mixture models to clustering. Firstly, they are not easily adapted to cope with discrete data. This is unfortunate because many real clustering problems involve both continuous and discrete variables. Secondly, they lead to models with large numbers of parameters: for example, if  $p = 8$  we need to estimate 36 parameters for even a common covariance matrix; many more if they must be estimated separately for each group. A third consideration is the common occurrence of missing values in multivariate data, particularly when the observations are on humans. A variant of MULTIMIX accommodates observations that are missing at random using the method of Little & Rubin (1987), but this is not described here.

Highly parameterized models can lead to difficulties in several ways. As discussed by McLachlan & Basford (1988 p. 11) the likelihood function of a mixture model can have singularities in a neighbourhood of which it is unbounded. Iterative methods for computing ML estimates are drawn towards these singularities from many starting values if the model is highly parameterized. It is also common to find many local maxima in such models. Even if we find the largest of the local maxima we often find that the likelihood is nearly constant in a low-dimensional set in which some of the parameters are functions of the others.

## 2.3. Mixtures of Discrete Distributions

Latent class analysis was developed by the mathematical sociologist Paul Lazarsfeld (Lazarsfeld & Henry, 1968) who was interested in making more precise the relationship between underlying or latent states that are not observable, and directly observable categorical variables indicating these states.

Latent class models can be described as follows: assume the population is made up of  $K$  groups or subpopulations  $G_1, \dots, G_K$  in proportions  $\pi_1, \dots, \pi_K$ . Let  $\mathbf{x}$  denote the vector

of responses on the  $p$  variables that we observe on each observation, where the  $j$ th variable can take on levels numbered from 1 to  $M_j$ . If the  $i$ th observation comes from subpopulation  $G_k$ , i.e.  $i \in G_k$ , then

$$f_k(\mathbf{x}_i; \boldsymbol{\theta}_k) = \Pr(\mathbf{X}_i = \mathbf{x}_i \mid i \in G_k) = \prod_{j=1}^p \Pr(X_{ij} = x_{ij} \mid i \in G_k)$$

where  $\boldsymbol{\theta}_k$  are the parameters of the distribution of the responses in the  $k$ th subpopulation; in this case, the set of probabilities:

$$\lambda_{kjm} = \Pr(X_{ij} = m \mid i \in G_k) \quad (m = 1, \dots, M_j, j = 1, \dots, p).$$

The overall probability function for the  $i$ th observation is a mixture of these conditional probability functions:

$$f(\mathbf{x}_i; \boldsymbol{\phi}) = \sum_{k=1}^K \pi_k f_k(\mathbf{x}_i; \boldsymbol{\theta}_k)$$

so that the latent class model is a finite mixture model. The parameter vector  $\boldsymbol{\phi}$  is made up of the  $\pi_k$  and the  $\lambda_{kjm}$  as the  $k$ ,  $j$  and  $m$  take on allowable values. Note that  $\sum_k \pi_k = 1$ ; and  $\sum_m \lambda_{kjm} = 1$  for any given  $j, k$ .

The original method of fitting these models, discussed at some length by Lazarsfeld and Henry for the case of binary variables, was to attempt to solve the system of equations given by equating the fitted cell probabilities to the observed cell proportions. The solution of these equations can be difficult and latent class analysis became much easier to use when Goodman (1974) introduced a new iterative algorithm for the ML fitting of latent class models. This algorithm is a special case of the general EM algorithm.

To use latent class analysis as a clustering method the probability  $\tau_{ik}$  that the  $i$ th observation comes from the  $k$ th group is first estimated by Bayes' rule from the estimated component distributions and the estimated proportions in each component. In fact these probabilities are also required in the course of the algorithm, although it is not until the algorithm has converged that we can use them for clustering. The versatility of latent class analysis as a clustering method was shown by Aitkin *et al.* (1981) who fitted two-class and three-class models to 38 binary variables describing how each of 468 teachers ran their classrooms, interpreting the classes as levels of a 'teaching style' factor in subsequent analyses. Pickering & Forbes (1984) used this method to study an even larger dataset. It consisted of clinical and diagnostic information about approximately 50 000 infant births. Eleven categorical variables each having from two to four levels were used to fit models having between one and six latent classes. The analysis was feasible because only about 600 distinct response profiles actually occurred in the data. Pickering & Forbes (1984) give references to other studies using latent class methods.

Most applications of latent class analysis remain within the social sciences where the method was developed. The ability to fit latent class models is one of the capabilities of  $\ell$ EM, a general program for fitting models to categorical data (Vermunt, 1997).

## 2.4. Everitt's Model for Ordinal Variables

Everitt (1988) proposed incorporating binary and ordinal variables into mixture models by means of 'threshold' parameters which divide the real line into regions corresponding to outcomes of the ordinal variable. Such threshold models have been widely used for ordinal data and a brief survey was given by Zhaorong, McGilchrist & Jorgensen (1992) where they

are used in a continuous latent variable model for the comparison of 20 ternary variables representing variants of microbiological test methods. These data could also have been analysed by latent class analysis which involves a discrete latent variable. Everitt & Mérette (1990) compared the clustering performance of Everitt's finite mixture method on four simulated datasets each having three continuous and two categorical variables, and on Fisher's iris data (Andrews & Herzberg, 1985 pp. 5–8) after two of the four variables had been categorized. They reported good performance of the mixture method compared with conventional hierarchical methods. There are severe practical limitations to the use of this method at present. Everitt proposed the use of standard optimization algorithms applied to the log-likelihood function. The computation of the log-likelihood function requires the numerical evaluation of a  $q$ -dimensional integral, where  $q$  is the number of categorical variables, and Everitt & Mérette (1990) considered no examples where  $q > 2$ . Their methods would be difficult to apply to the highly multivariate datasets to which cluster analysis has traditionally been applied. For these reasons, MULTIMIX makes no special provision for ordinal variables. Depending on the circumstances, it is usually acceptable to treat them either as categorical or as continuous.

### 3. The MULTIMIX Model Family

We now describe in detail our initial approach to mixture model clustering. We expect the data to be in the form of an  $n \times p$  matrix of observations by variables that we regard as a random sample from the distribution  $f(\mathbf{x}) = \sum \pi_k f_k(\mathbf{x})$ , itself a finite mixture of  $K$  component distributions  $f_k$  in proportions  $\pi_k \geq 0$  satisfying  $\sum \pi_k = 1$ . We suppose that the vector of variables  $\mathbf{x} = (x_1, \dots, x_p)'$  has been partitioned into  $(\check{\mathbf{x}}_1' | \dots | \check{\mathbf{x}}_L')'$ . We consider component distributions that factorize  $f_k(\mathbf{x}) = \prod_{\ell} f_{k\ell}(\check{\mathbf{x}}_{\ell})$ , conformably with this partition. This is a weak form of 'local independence': within each of the  $K$  subpopulations the variables in the subvector  $\check{\mathbf{x}}_{\ell}$  are independent of the variables in  $\check{\mathbf{x}}_{\ell'}$  for  $\ell \neq \ell'$ . True local independence is the independence of each  $x_j$  within subpopulations. We can write the model for the  $i$ th observation as

$$f(\mathbf{x}_i; \boldsymbol{\phi}) = \sum_{k=1}^K \pi_k \prod_{\ell=1}^L f_{k\ell}(\check{\mathbf{x}}_{i\ell}; \boldsymbol{\theta}_{k\ell}), \quad (2.1)$$

where  $\boldsymbol{\theta}_{k\ell}$  consists of the parameters of the distribution  $f_{k\ell}$ , and the  $\pi_k$  are the mixing proportions. This formulation includes the motivating examples of latent class analysis (Aitkin *et al.*, 1981) and mixtures of multivariate normals (McLachlan & Basford, 1988). With one exception to be described later, subvectors are usually formed with vectors of the same type, categorical or continuous. When a subvector contains only a single variable, that variable is independent of all other variables within each subpopulation.

It is convenient to assume forms for the  $f_{k\ell}$ , and hence for the  $f_k$ , that belong to the exponential family. The model is then well suited for ML estimation of its parameters by the EM algorithm of Dempster *et al.* (1977). This approach is followed in MULTIMIX with the following distributions for the  $\check{\mathbf{x}}_{k\ell}$ :

- (a) *discrete distribution*. Here  $\check{\mathbf{x}}_{\ell} = u_{\ell}$  is a one-dimensional discrete random variable taking values  $1, \dots, M_j$  with probabilities  $\lambda_{k\ell 1}, \dots, \lambda_{k\ell M_j}$ .
- (b) *multivariate normal*. Here  $\check{\mathbf{x}}_{\ell} = \mathbf{v}_{\ell}$  is a  $p_{\ell}$ -dimensional vector of continuous random variables with the  $N_{p_{\ell}}(\mu_{k\ell}, \Sigma_{k\ell})$  distribution.



(c) *location model*. Here  $\tilde{\mathbf{x}}_\ell = (u_\ell, \mathbf{v}_\ell)$  is a  $(1 + p_\ell)$ -dimensional vector of random variables with one discrete variable,  $x_j$ , and  $p_\ell$  continuous variables as elements. The discrete random variable takes values  $1, \dots, M_j$  with probabilities  $\lambda_{k\ell 1}, \dots, \lambda_{k\ell M_j}$ . Conditional on the discrete variable taking value  $m$  the  $p_\ell$  continuous random variables have the multivariate normal distribution  $N_{p_\ell}(\nu_{mk\ell}, \Xi_{k\ell})$ ; see Krzanowski (1983) for details.

If all variables are of continuous type, then  $f(\mathbf{x}) = \sum \pi_k f_k(\mathbf{x})$  is a mixture of multivariate normal distributions. The way in which the set of variables is partitioned into subvectors determines the form of the matrix of covariance parameters in each  $f_k$ . The form is block-diagonal with a square block corresponding to each subvector. Extreme cases are the fully unstructured covariance matrix case and the diagonal covariance matrix case. Unstructured covariance matrices introduce many parameters into the model and hence should be avoided as far as possible. A reasonable strategy for fitting a mixture of multivariate normals for clustering purposes would be to begin with the local independence case (diagonal covariance matrices) and then to estimate the model parameters, assign observations to clusters and then study the within-cluster correlation matrices. Variables that are highly correlated in some of the clusters could be grouped into a subvector and the whole process repeated with the model so modified.

If all the variables are of discrete type, the model is the usual latent class model. In principle, local independence could fail in this situation as well, although this is not often checked for. If strong within-cluster associations between two discrete variables are detected after a preliminary clustering then the two variables may be combined into a single discrete variable with a level for each cell of the two-way table (or fewer, if some cells are pooled).

The location model for a subvector in the partition is introduced in the general MULTIMIX model to cope with the possibility of within-cluster associations between a discrete variable and several continuous variables. We do not expect this facility to be needed very often in practice. Other types of variable models are available in other mixture modelling programs to be discussed below. It is no problem, in principle, to add new types of attribute model from within the exponential family to extend the MULTIMIX model. As the model has been described, it is a mixture of  $K$  distributions, each of which can be seen to belong to the exponential family. It is therefore well suited for ML estimation of its parameters by the EM algorithm of Dempster *et al.* (1977), and the FORTRAN 77 program MULTIMIX was written by Hunt (1996) to do this.

The ‘complete data’, in EM terminology, consists of the  $n \times p$  array of observed data  $\{x_{ij}\}$  and the conceptual  $n \times K$  array  $\{z_{ik}\}$  of class membership indicators. The indicator vectors  $\mathbf{z}_1, \dots, \mathbf{z}_n$  are independently and identically distributed according to a multinomial distribution generated by one draw on a population made up of  $K$  categories in proportions  $\pi_1, \dots, \pi_K$ .

The complete-data specification treats the  $\mathbf{z}_i$  as known leading to the log-likelihood

$$\begin{aligned} L_C(\boldsymbol{\phi}) &= \log \left( \prod_{i=1}^n \prod_{k=1}^K \left[ \pi_k^{z_{ik}} \left\{ \prod_{\ell=1}^L f_{k\ell}(\mathbf{x}_i; \boldsymbol{\theta}_{k\ell}) \right\}^{z_{ik}} \right] \right) \\ &= \sum_{i=1}^n \sum_{k=1}^K \left\{ z_{ik} \log \pi_k + z_{ik} \sum_{\ell=1}^L \log f_{k\ell}(\mathbf{x}_i; \boldsymbol{\theta}_{k\ell}) \right\} \\ &= \sum_{i=1}^n \sum_{k=1}^K z_{ik} \log \pi_k + \sum_{k=1}^K \ell_k(\boldsymbol{\theta}_k) \end{aligned}$$

$$\text{where} \quad \ell_k(\theta_k) = \sum_{i=1}^n \left\{ z_{ik} \sum_{\ell=1}^L \log f_{k\ell}(x_i; \theta_{k\ell}) \right\} = \sum_{\ell=1}^L \sum_{i=1}^n z_{ik} \log f_{k\ell}(x_i; \theta_{k\ell}).$$

Maximizing the complete-data log-likelihood  $L_C(\phi)$  is equivalent to maximizing  $\ell_k(\theta_k)$  separately for each subvector in the partition. By substituting the appropriate density for the  $f_{k\ell}$ , Hunt (1996) deduces that the complete-data sufficient statistics for the model are the following:

1. for each class  $G_k$ :

$$\sum_i z_{ik};$$

2. for each class  $G_k$ , each categorical variable  $x_j$ , and each value  $m$  of  $x_j$ :

$$\sum_i z_{ik} \delta_{ijm}, \quad \text{where } \delta_{ijm} = \begin{cases} 1 & \text{if } x_{ij} = m, \\ 0 & \text{otherwise;} \end{cases}$$

3. (a) for each class  $G_k$  and each continuous variable  $x_j$  belonging to a multivariate normal subvector:

$$\sum_i z_{ik} x_{ij} \quad \text{and} \quad \sum_i z_{ik} x_{ij}^2;$$

- (b) for each class  $G_k$  and each pair of continuous variables  $x_j$  and  $x_{j'}$ ,  $j < j'$ , belonging to the same multivariate normal subvector:

$$\sum_i z_{ik} x_{ij} x_{ij'};$$

4. (a) for each class  $G_k$ , each continuous variable  $x_j$  belonging to a location model subvector indexed by  $\ell$  and each value  $m$  of the categorical variable  $u_\ell$ :

$$\sum_i z_{ik} w_{i\ell m} x_{ij} \quad \text{and} \quad \sum_i z_{ik} w_{i\ell m} x_{ij}^2;$$

- (b) for each class  $G_k$ , each pair of continuous variables  $x_j$  and  $x_{j'}$ ,  $j < j'$ , belonging to the location model subvector indexed by  $\ell$  and each value  $m$  of the categorical variable  $u_\ell$ :

$$\sum_i z_{ik} w_{i\ell m} x_{ij} x_{ij'} \quad \text{where } w_{i\ell m} = \begin{cases} 1 & \text{if } u_\ell = m, \\ 0 & \text{otherwise.} \end{cases}$$

The EM iteration alternates between two calculations, the E-step and the M-step. Beginning at a current value for  $\phi$ , say  $\phi^{(p)}$  the vector of all unknown parameters, the E-step requires the calculation of  $Q(\phi, \phi^{(p)}) = E\{L_C(\phi) \mid X; \phi^{(p)}\}$ , the expectation of the complete-data log-likelihood, conditional on the observed data and the current value of the parameters. Because the complete-data sufficient statistics are linear in the unobserved  $z_{ik}$  we can calculate  $Q(\phi, \phi^{(p)})$  from  $L_C(\phi)$  by replacing  $z_{ik}$  with

$$\hat{\tau}_{ik} = E(z_{ik} \mid \mathbf{x}_i; \phi^{(p)}) = \frac{\pi_k^{(p)} f_k(\mathbf{x}_i, \theta_k^{(p)})}{\sum_{k=1}^K \pi_k^{(p)} f_k(\mathbf{x}_i, \theta_k^{(p)})}$$

in  $L_C(\phi)$ . That is,  $z_{ik}$  is replaced by the estimate of the posterior probability  $\tau_{ik}$  that observation  $i$  belongs to group  $G_k$ .

At the M-step  $\phi^{(p+1)}$  is chosen to be a value of  $\phi$  which maximizes  $Q(\phi, \phi^{(p)})$  with respect to its first argument. For the MULTIMIX model the elements of  $\phi^{(p+1)}$  are given by

$$\begin{aligned}\hat{\pi}_k &= \frac{1}{n} \sum_{i=1}^n \hat{\tau}_{ik}, & \hat{\lambda}_{k\ell m} &= \frac{1}{n\hat{\pi}_k} \sum_{i, u_{i\ell}=m} \hat{\tau}_{ik}, \\ \hat{\mu}_{k\ell} &= \frac{1}{n\hat{\pi}_k} \sum_{i=1}^n \hat{\tau}_{ik} \mathbf{v}_{i\ell}, & \hat{\Sigma}_{k\ell} &= \frac{1}{n\hat{\pi}_k} \sum_{i=1}^n \hat{\tau}_{ik} (\mathbf{v}_{i\ell} - \hat{\mu}_{k\ell})(\mathbf{v}_{i\ell} - \hat{\mu}_{k\ell})', \\ \hat{\nu}_{k\ell m} &= \frac{1}{n\hat{\pi}_k} \sum_{i, u_{i\ell}=m} \hat{\tau}_{ik} \mathbf{v}_{i\ell}, & \hat{\Xi}_{k\ell} &= \frac{1}{n\hat{\pi}_k} \sum_{i, u_{i\ell}=m} \hat{\tau}_{ik} (\mathbf{v}_{i\ell} - \hat{\nu}_{k\ell m})(\mathbf{v}_{i\ell} - \hat{\nu}_{k\ell m})',\end{aligned}$$

for  $k = 1, \dots, K$  and  $\ell = 1, \dots, L$ . Note that the level probabilities  $\lambda_{k\ell m}$  for the categorical variables are calculated in the same way, irrespective of whether or not the discrete variable  $u_\ell$  belongs to a location model subvector.

The current version of the program uses a convergence criterion to cease iterating when the difference in log-likelihoods at iteration  $t$  and iteration  $t - 10$  is less than 0.0000001. The iteration may be started either from an initial classification or from an initial set of parameter estimates. As the number of parameters is quite large it is usually more convenient to begin with a classification.

#### 4. Example: Byar Prostate Cancer Data

We consider the clustering of cases on the basis of pre-trial covariates alone for the prostate cancer clinical trial data of Byar & Green (1980) reproduced in Andrews & Herzberg (1985 pp. 261–274). This dataset was obtained from a randomized clinical trial comparing four treatments for 506 patients with prostatic cancer grouped on clinical criteria into stages 3 and 4 of the disease. As reported by Byar & Green, stage 3 represents local extension of the disease without evidence of distant metastasis, while stage 4 represents distant metastasis as evidenced by elevated acid phosphatase, X-ray evidence, or both. We assess the clusters obtained by MULTIMIX in comparison with the clinical stages, and also consider the trial outcomes for patients in different clusters. The treatments consisted of oestrogen therapy at differing rates. Daily pills containing 0.0 (placebo), 0.2, 1.0, and 5.0 mg of diethylstilbestrol were administered in the four treatments. Byar and Green noted little difference between the effects of the first two treatments, nor between the effects of the last two treatments, so we call patients in either of the first two treatments ‘Untreated’ and in either of the last two treatments ‘Treated’.

There are twelve pre-trial covariates (Table 1) measured on each patient, seven may be taken to be continuous, four to be discrete, and one variable (SG) is an index nearly all of whose values lie between 7 and 15, and which could be considered either discrete or continuous. We treat SG as a continuous variable. A preliminary inspection of the data shows that the size of the primary tumour (SZ) and serum prostatic acid phosphatase (AP) are both skewed variables. These variables have therefore been transformed, SZ by a square root transformation, and AP using a logarithmic transformation, to make their distributions more symmetric. Observations that are missing values in any of the twelve pretreatment covariates are omitted from further analysis, leaving 475 out of the original 506 observations available. In fact, several of the

TABLE 1  
*Pretreatment covariates*

Covariate	Abbreviation	Number of levels (if categorical)
Age	Age	
Weight index	WtI	
Performance rating	PF	4
Cardiovascular disease history	HX	2
Systolic blood pressure	SBP	
Diastolic blood pressure	DBP	
Electrocardiogram code	EKG	7
Serum haemoglobin	HG	
Size of primary tumour	SZ	
Index of tumour stage and histologic grade	SG	
Serum prostatic acid phosphatase	AP	
Bone metastases	BM	2

TABLE 2  
*Models and fits*

Model	Variable groups	No. parameters (#P)	Log-likelihood +11386.265	$\frac{\Delta \log L}{\Delta \#P}$
[LInd]	—	55	0.000	58.77
[BPr]	{ SBP, DBP }	57	117.542	5.31
[3,2]	{ BM, WtI, HG }, { SBP, DBP }	63	149.419	1.65
[5]	{ BM, WtI, HG, SBP, DBP }	75	169.163	1.31
[9]	Complement of { PF, HX, EKG }	127	237.092	

analyses to be described were also carried out using a version of the program which allows for missing observations, treating them as missing at random in the sense of Little & Rubin (1987). There was little variation from the results using only the complete observations.

We consider the fitting of two-class models ( $K = 2$ ). The simplest model is the model [LInd] of complete local independence in which the component densities take the form

$$f_k(\mathbf{x}_i; \boldsymbol{\theta}_k) = \prod_{\ell=1}^{12} f_{k\ell}(\check{\mathbf{x}}_{i\ell}; \boldsymbol{\theta}_{k\ell}),$$

where  $\boldsymbol{\theta}_{k\ell}$  is the parameter vector for subvector  $\ell$ , group  $k$ ,  $k = 1, 2$ . Note that  $f_{k\ell}(\check{\mathbf{x}}_{i\ell}; \boldsymbol{\theta}_{k\ell})$  is  $N(\mu_{k\ell}, \sigma_{k\ell}^2)$  for each of the eight continuous variables, and  $D(\lambda_{k\ell 1}, \dots, \lambda_{k\ell m_\ell})$  for each of the four categorical variables.

The fitting strategy used was a form of forward selection of covariances, beginning with [LInd] and progressively adding local associations to the model by taking coarser and coarser partitions of the set of 12 variables. The modifications to the current model were determined by examining correlations, scatterplots and two-way tables within each of the two clusters formed by allocating each observation according to the current model. Table 2 summarizes the results of this fitting process and a description of some of the steps follows.

When the data had been grouped into two classes following the fitting of [LInd], correlations between SBP and DBP of about 0.62 were observed within both of the classes, and these

TABLE 3  
*Clusters and outcomes for Treated and Untreated patients*

Patient group		Outcome			
		Alive	Prostate Death	Cardio Death	Other Death
Untreated patients					
Cluster 1	Stage 3	39	18	37	33
	Stage 4	3	4	3	3
Cluster 2	Stage 3	1	4	2	3
	Stage 4	14	49	18	6
Treated patients					
Cluster 1	Stage 3	50	3	52	20
	Stage 4	4	0	1	3
Cluster 2	Stage 3	1	6	3	1
	Stage 4	25	37	22	10

appeared to be the strongest associations. The fact that one would expect such a correlation within any naturally formed group of patients made it compelling to fit a model [BPr] in which SBP and DBP had a bivariate normal distribution within groups. The partition of the variables for this model placed these two variables together in a subvector, the remaining subvectors being singletons. Thus [BPr] contains two more parameters than [LInd].

The next group of variables chosen was the triple {BM, WtI, HG}, giving a location model factor to the mixture densities, because BM is dichotomous while WtI and HG are continuous. The resultant model is denoted by [3,2], referring to the size of these variable groups. Six extra parameters are introduced in this change: there are four new mean parameters, because the fitted means of WtI and HG are now specific to each level of BM within each subpopulation, and two new covariance parameters. Model [5] combines these two variable groups at a cost of introducing 12 new parameters, and Model [9] has one large variable group combining BM with all eight continuous variables. Table 2 also includes the log-likelihoods obtained. In the case of [LInd], [BPr], [3,2] and [5] these log-likelihoods were obtained from several initial configurations including random groupings of the observations; however, [9] proved to be sensitive to the choice of starting configuration and the greatest log-likelihood over four runs is shown for this model. Convergence was usually obtained after 60 to 70 iterations although one run for model [9] reached 200 iterations without converging.

There was little difference between the group allocations determined by [LInd], [BPr], [3,2] and [5], with the allocation of only four patients out of 475 changing between these models. Model [9] allocations were sensitive to the initial classification and did not agree so closely with each other nor with the classifications of the more parsimonious models. Comparing the [BPr] allocation with the clinical grouping into stages 3 and 4 of the disease we find one cluster with 252 stage 3 and 21 stage 4 patients and the other cluster with 21 stage 3 and 181 stage 4 patients.

It is of interest to examine the post-trial survival status of patients in the four stage-cluster combinations, which have been arrived at using pre-trial information only. Table 3 presents this information for the [BPr] model. While model cluster 1 and clinical stage 3 are associated with a better chance of survival, the patterns of outcomes for the 42 patients whose model and clinical classifications conflict suggest that the model classifications are better indicators of prognosis than the clinical criteria used. This is especially noticeable among the Treated patients.

Hunt (1996) analyses this dataset in more detail, also fitting three-class and four-class models yielding classifications with distinctively different outcome patterns, suggesting that the models were detecting real features of the population. She also develops the methods of Little & Rubin (1987) for use with the model of this paper and applies this to the complete dataset of 506 prostate cancer patients.

## 5. Model Comparison Tests

### 5.1. Number of Components in the Mixture

McLachlan & Basford (1988 Sect. 1.10) considers the question of testing for the number of components in a mixture. The problem is difficult because although a model with  $K_1$  components is nested within a model with  $K_2 > K_1$  components, the usual regularity conditions are not met. These conditions are required to conclude that if  $\lambda$  is the likelihood ratio,  $-2 \log \lambda$  is asymptotically distributed as  $\chi^2$  with degrees of freedom equal to the difference in the number of parameters in the two models. In fact, the asymptotic distribution may depend on the true values of the parameters of the component distributions, so there is no general result. McLachlan & Basford discuss a number of bootstrap approaches to the problem. Feng & McCulloch (1992, 1994, 1996) have studied several aspects of this problem. In their most recent paper they recommend a bootstrap procedure. Bootstrap procedures would be very costly to apply to the clustering of datasets with many observations on many variables. Wolfe (1971) investigated the distribution of  $-2 \log \lambda$  when comparing nested mixtures of multivariate normal distributions and recommended treating the distribution as  $\chi^2$ , but with double the nominal degrees of freedom. Banfield & Raftery (1993) developed a Bayesian *approximate weight of evidence* criterion as a guide for choosing the number of components in the mixture. Wallace & Dowe (1998) use a 'minimum message length' criterion as a basis for their parameter estimation. This method unifies model selection and parameter estimation and leads to a choice for the number of components.

A specific example provides an illustration of how it may not be realistic to expect to choose a value for  $K$  on sample evidence alone. Consider the problem of estimating growth and age structure in a stock of fish from length–frequency data. Suppose that  $k = 1, \dots, K$  indexes  $K$  subpopulations  $\Pi_k$  (age classes) of fish and that the fish in  $\Pi_k$  all have age  $t_k$  years, where  $t_k = t_1 + k - 1$ . Let  $\pi_k$  be the proportion of the population in  $\Pi_k$  and  $\mu_k$  and  $\sigma_k$  be the mean and standard deviation of  $\Pi_k$ . Schnute & Fournier (1980) discuss the ML fitting of a model of this kind where the parameters  $\mu_k$  and  $\sigma_k$  are modelled parametrically as functions of  $t_k$ . In the fisheries application discussed by Schnute and Fournier, the  $\mu_k$  and  $\sigma_k$  tend to limiting values as  $k$  increases and the  $\pi_k$  vary because of annual recruitment variations but tend to diminish geometrically because of cumulative mortality, both natural and through fishing. Thus as  $k$  increases the  $\Pi_k$  become closer together but are represented in the sample by smaller and smaller proportions. This kind of situation seems very natural but would appear to resist any form of statistical inference for the value of  $K$ . In view of the complexities of this question it seems best to regard the number of components  $K$  as a choice to be made by the modeller, in much the same way as a functional form for a distribution is selected arbitrarily. This does not mean that model comparison statistics of the  $-2 \log \lambda$  kind cannot be used heuristically. There remains the possibility that a small number of observations from unmodelled components will upset the fit of the model to the bulk of the data. Jorgensen (1990) discussed a number of diagnostic statistics that may be used to detect these points.

## 5.2. Number of Within-cluster Associations

In contrast to the determination of the number of classes  $K$ , the standard likelihood ratio tests for two nested models based on an approximate distribution for  $-2 \log \lambda$  of  $\chi^2$  with degrees of freedom equal to the difference in the number of parameters in the two models are not likely to mislead. The most troublesome regularity condition requiring checking is that the third order derivatives of  $\log f$  with respect to the parameters are bounded (in a neighbourhood of the true parameter vector) by functions of the data with finite expectation (Lehmann, 1983 p.429). This can be shown to be the case, roughly, if the partial derivatives of the component densities  $f_k$ , with respect to all parameters and up to third order, are not too large in comparison with the mixture density  $f$ . Even for the simple case of a mixture of two bivariate normals, a model with 11 parameters, we have  $6^3 + 6^3$  third order derivatives of component densities to check, although many of these coincide. Checking some of these by hand suggests that all will be well unless a true variance parameter is zero or a correlation is  $\pm 1$  — that is, unless a component density is degenerate. It is also necessary that neither proportion be close to zero. A natural conjecture in the case of a finite mixture of multivariate normals is that the regularity conditions are satisfied as long as the smallest eigenvalue of the true variance–covariance matrix for each group is not close to zero. In practical terms, the suggestion is that when a ‘Reduced’ model is being compared with a ‘Full’ model (having the same number of components, but with extra association parameters), we may base a model comparison test on the assumption that  $-2 \log \lambda$  is approximately distributed as  $\chi^2$  with degrees of freedom equal to the difference in the number of parameters in the two models when the Reduced model is operating, *unless* the fit obtained under the Reduced model has any degeneracies either in the number of components or in the form of any of the components. Dempster (1972) and Wermuth (1976a, b) considered similar model selection problems, but for the case of a single multivariate normal component rather than a mixture of these. Wermuth also considered loglinear contingency table models. These authors parameterized the multivariate normal using the inverse of the covariance matrix, called the concentration matrix, rather than the covariance matrix. They tested for the vanishing of a set of elements of the concentration matrix, which is equivalent to the vanishing of the corresponding set of partial correlations. We restrict ourselves to models with block-diagonal covariance structure, and the inverses may be calculated block by block, so tests involving the splitting or combining of blocks may be formulated in terms of either the covariances or the concentrations.

## 6. Other Programs for Mixture Model Clustering

### 6.1. AUTOCLASS

AUTOCLASS (Cheeseman & Stutz, 1996) is a Bayesian clustering program. The models fitted by AUTOCLASS are very similar to those fitted by MULTIMIX, although both programs were developed independently. Two obvious differences are that (i) AUTOCLASS has automated the process of model selection as well as that of parameter estimation but MULTIMIX leaves model-specification to the user; and (ii) AUTOCLASS uses maximum posterior (MAP) estimation in place of ML estimation.

In fact the first is the more crucial difference, because the EM algorithm that is the basis of both programs accommodates both ML and MAP estimation. AUTOCLASS compares different models by calculating an approximation to the marginal density of the observed data after the model parameters have been integrated out. In usual EM language

the approximation used is analogous to taking observed data likelihood to be proportional to complete data likelihood with the constant of proportionality to be evaluated at the ML estimates.

The models currently available in AUTOCLASS for attributes within a component are as follows. Categorical attributes are modelled by general discrete distributions (multi-category Bernoulli) as in MULTIMIX. Continuous attributes may be taken to have uniform or normal distributions, possibly after transformation. Poisson distributions are available for count attributes. Cheeseman & Stutz (1996) report that von Mises–Fisher distributions for circular and spherical attributes are under development. At present it appears that AUTOCLASS does not offer facilities for modelling within-cluster dependencies; that is, all models assume within-cluster independence of attributes. Missing values are treated as a special kind of value in some attribute models, but there has been no implementation of the Little & Rubin (1987) method for data missing at random.

Cheeseman and Stutz claim that the AUTOCLASS method of model comparison introduces an ‘Occam factor’ which penalizes overfitting. However, Edwards & Dowe (1998) describe the minimum message length (MML) fitting of a model that combines a continuous latent factor with a number of classes to a set of 5425 infrared spectra from astronomical point sources. Edwards & Dowe find 12 classes where AUTOCLASS found 77 (Goebel *et al.*, 1989). It is not clear whether the difference in the number of classes in the fit is due to the explicit penalty on overfitting built into the MML criterion or whether it is the introduction of the continuous factor which is responsible.

## 6.2. SNOB

SNOB (Wallace & Dowe, 1998) is a clustering program whose development by workers began in the late sixties (Wallace & Boulton, 1968). SNOB could be said to be class-conscious. SNOB\* is a mixture model similar in structure to AUTOCLASS and offering local independence models based on discrete, normal, Poisson and von Mises distributions. In fact, SNOB is the older program. A novel feature of SNOB is that inference is by the principle of minimum message length (Wallace & Freeman, 1987). This form of inference takes discrete variables as fundamental and seeks to minimize the negative logarithm of the probability of the model and parameter values, plus the negative logarithm of the probability of the data given the model and parameter values. A continuous analogue of this estimation principle is similar to MAP estimation but introduces an additional factor,  $F(\theta)^{-1/2}$ , to the prior, where  $F(\theta)$  is the determinant of the Fisher information matrix at the parameter vector  $\theta$ .

In contrast to MULTIMIX, where the user must specify the number of classes, SNOB selects the number of classes automatically using the MML criterion. Thus the MML criterion is used for all aspects of model selection and parameter estimation in the SNOB approach.

## 6.3. MCLUST

Banfield & Raftery (1993) developed the classification likelihood approach of Scott & Symons (1971) further to introduce a controlled amount of flexibility to criterion-based cluster analysis for continuous data. This approach suffers from the disadvantages mentioned in Section 1.2, but it does lead to similar optimization problems to those faced in traditional cluster analysis, and hence the model may be fitted by algorithms similar to those used to solve those

---

\* SNOB has a home page at <http://www.cs.monash.edu.au/~dld/Snob.html>.



problems. Wallace & Dowe (1998) point out that in the case of a substantially overlapping pair of normal distributions having equal abundance and common standard deviation, this kind of estimation is likely to overestimate the difference in means and underestimate the standard deviation.

Banfield and Raftery characterize the dispersion matrices of multivariate normal clusters by their *orientation*, *size*, and *shape*. They mainly consider models where the shape is the same in each component of the mixture, but orientation and size are permitted to vary. They also consider an approach to robustifying cluster analysis by allowing a very dispersed 'noise' component in addition to the multivariate normal components.

A FORTRAN program called MCLUST has been written by Fraley to fit these models and others.<sup>†</sup> Although criterion-based, rather than being based on a distance matrix, MCLUST is written to proceed initially as an agglomerative hierarchical program. However, once the number of clusters has been determined by the user, MCLUST can proceed by reallocating points to seek a minimum of the criterion in a fashion similar to the *k*-means algorithm of Hartigan (1975). In recent versions of S-PLUS, MCLUST now forms the core of the clustering functions provided.

## 7. The Place of MULTIMIX in Mixture Modelling

The brief survey of other related programs helps to clarify the role of MULTIMIX as a mixture modelling tool. In contrast to SNOB and AUTOCLASS it automates only parameter estimation, leaving model selection to the control of the user. It appears to be unique in offering a ML approach to a class of models extending mixtures of multivariate normals and latent class models (although it is possible that AUTOCLASS and SNOB might be coaxed into producing similar output for at least some models by appropriate prior specification and the switching off of their model search facilities).

A natural further development for MULTIMIX would be to introduce new types of attribute distribution such as the Poisson and circular von Mises distributions. To the extent that robust estimation is appropriate for a particular dataset it seems that it would be better to add a very small proportion of a highly dispersed component to the mixture than to follow Banfield & Raftery (1993) in modifying the likelihood criterion to gain robustness.

There are no present plans to automate model selection in MULTIMIX, but it must be acknowledged that more needs to be done in the way of graphical diagnostic output to assist the user with the refinement of the models. Eventually some form of automation of model selection will be necessary if MULTIMIX is to be used on extremely large datasets, but we would feel happier about adopting any proposal for model selection if we could compare it with human-driven procedures over a range of datasets.

The availability of the four programs AUTOCLASS, MCLUST, MULTIMIX and SNOB offering similar ranges of models but using different inferential principles provides an opportunity to learn more about the strengths and weaknesses of these principles in the practical data analysis context of large multivariate datasets. Currently MULTIMIX is available as FORTRAN 77 source code from the authors<sup>‡</sup> as are some documentation, datasets and auxiliary programs.

<sup>†</sup> MCLUST is available from StatLib either as a FORTRAN program or as an S-PLUS function.

<sup>‡</sup> URL <ftp://ftp.math.waikato.ac.nz/pub/maj/>

## References

- AITKIN, M., ANDERSON, D. & HINDE, J. (1981). Statistical modelling of data on teaching styles. *J. Roy. Statist. Soc. Ser. A* **144**, 419–461.
- ANDERBERG, M.R.C. (1973). *Cluster Analysis for Applications*. New York: Academic Press.
- ANDREWS, D.A. & HERZBERG, A.M. (1985). *Data: a Collection of Problems from Many Fields for the Student and Research Worker*. Springer Series in Statistics. New York: Springer-Verlag.
- BANFIELD, J.D. & RAFTERY, A.E. (1993). Model-based Gaussian and non-Gaussian clustering. *Biometrics* **49**, 803–821.
- BINDER, D.A. (1978). Bayesian cluster analysis. *Biometrika* **65**, 31–38.
- BRYANT, P. & WILLIAMSON, J.A. (1978). Asymptotic behaviour of classification maximum likelihood estimates. *Biometrika* **65**, 273–281.
- BUTLER, R.W. (1986). Predictive likelihood inference with applications (with discussion). *J. Roy. Statist. Soc. Ser. B* **48**, 1–38.
- BYAR, D.P. & GREEN, S.B. (1980). The choice of treatment for cancer patients based on covariate information: application to prostate cancer. *Bull. Cancer (Paris)* **67**, 477–490.
- CHEESEMAN, P. & STUTZ, J. (1996). Bayesian classification (AutoClass): theory and results. In *Advances in Knowledge Discovery and Data Mining*, eds U.M. Fayyad, G. Piatetsky-Shapiro, P. Smyth & R. Uthurusamy, 153–180. Cambridge, MA: AAAI Press/MIT Press.
- CORMACK, R.M. (1971). A review of classification (with discussion). *J. Roy. Statist. Soc. Ser. B* **134**, 321–367.
- DAY, N.E. (1969). Estimating the components of a mixture of normal distributions. *Biometrika* **56**, 463–474.
- DEMPSTER, A.P. (1972). Covariance selection. *Biometrics* **28**, 157–175.
- , LAIRD, N.M. & RUBIN, D.B. (1977). Maximum likelihood from incomplete data via the EM algorithm (with discussion). *J. Roy. Statist. Soc. Ser. B* **39**, 1–38.
- EDWARDS, R.T. & DOWE, D.L. (1998). Single factor analysis in MML mixture modelling. In *Research and Development in Knowledge Discovery and Data Mining*, eds X. Wu, R. Kotagiri & K.B. Korb, Vol. 1394 of *Lecture Notes in Artificial Intelligence*, 96–109. Berlin: Springer.
- EVERITT, B.S. (1980). *Cluster Analysis*. 2nd edn. New York: Halsted.
- (1985). Mixture distributions. In *Encyclopedia of Statistical Sciences*, eds S. Kotz, N.L. Johnson & C.B. Read, Vol. 5, 559–569. New York: Wiley.
- (1988). A finite mixture model for the clustering of mixed-mode data. *Statist. Probab. Lett.* **6**, 305–309.
- & HAND, D.J. (1981). *Finite Mixture of Distributions*. London: Chapman and Hall.
- & MÉRETTE, C. (1990). The clustering of mixed-mode data: a comparison of possible approaches. *J. Appl. Statist.* **17**, 283–297.
- FENG, Z.D. & McCULLOCH, C.E. (1992). Statistical inference using maximum likelihood estimation and the generalized likelihood ratio when the true parameter is on the boundary of the parameter space. *Statist. Probab. Lett.* **13**, 325–332.
- & — (1994). On the likelihood ratio test statistics for the number of components in a normal mixture with unequal variance. *Biometrics* **50**, 1158–1162.
- & — (1996). Using bootstrap likelihood ratios in finite mixture models. *J. Roy. Statist. Soc. Ser. B* **58**, 609–617.
- FRIEDMAN, H.P. & RUBIN, J. (1967). On some invariant criteria for grouping data. *J. Amer. Statist. Assoc.* **62**, 1152–1178.
- GOEBEL, J., VOLK, K., WALKER, H., GERBAULT, F., CHEESEMAN, P., SELF, M., STUTZ, J. & TAYLOR, W. (1989). A Bayesian classification of the IRAS LRS atlas. *Astronom. Astrophys.* **225**, L5–L8.
- GOODMAN, L.A. (1974). Exploratory latent structure analysis using both identifiable and unidentifiable models. *Biometrika* **61**, 215–231.
- GORDON, A.D. (1981). *Classification*. In *Monographs on Statistics and Applied Probability*, Vol. 16. London: Chapman and Hall.
- GOWER, J.C. (1971). A general coefficient of similarity and some of its properties. *Biometrics* **27**, 857–871.
- & ROSS, G.J.S. (1969). Minimum spanning trees and single linkage cluster analysis. *Appl. Statist.* **18**, 54–64.
- HAND, D.J. (1981). *Discrimination and Classification*. London: Wiley.
- HARTIGAN, J.A. (1975). *Clustering Algorithms*. New York: Wiley.
- HASSELBLAD, V. (1966). Estimation of parameters for a mixture of normal distributions. *Technometrics* **8**, 431–444.
- (1969). Estimation of finite mixtures of distributions from the exponential family. *J. Amer. Statist. Assoc.* **64**, 1459–1471.

- HUNT, L.A. (1996). *Clustering using Finite Mixture Models*. PhD thesis, University of Waikato.
- JAMES, M. (1985). *Classification Algorithms*. London: Collins.
- JARDINE, N. & SIBSON, R. (1971). *Mathematical Taxonomy*. London: Wiley.
- JEFFREYS, H. (1932). An alternative to the rejection of observations. *Proc. Roy. Soc. London Ser. A* **137**, 78–87.
- JORGENSEN, M.A. (1990). Influence based diagnostics for finite mixture models. *Biometrics* **46**, 1047–1058.
- KRZANOWSKI, W.J. (1983). Distance between populations using mixed continuous and categorical variables. *Biometrika* **70**, 235–243.
- LANCE, G.N. & WILLIAMS, W.T. (1967). A general theory of classificatory sorting strategies. I. Hierarchical systems. *Comput. J.* **9**, 373–380.
- LAZARSFELD, P.F. & HENRY, N.W. (1968). *Latent Structure Analysis*. Boston: Houghton Mifflin.
- LEHMANN, E.H. (1983). *Theory of Point Estimation*. New York: Wiley.
- LITTLE, R.J.A. & RUBIN, D.B. (1987). *Statistical Analysis with Missing Data*. New York: Wiley.
- MARDIA, K.V., KENT, J.T. & BIBBY, J.M. (1979). *Multivariate Analysis*. London: Academic Press.
- MARRIOT, F.H.C. (1975). Separating mixtures of normal distributions. *Biometrics* **31**, 767–769.
- McLACHLAN, G.J. (1992). *Discriminant Analysis and Statistical Pattern Recognition*. New York: Wiley.
- & BASFORD, K.E. (1988). *Mixture Models: Inference and Applications to Clustering*. New York: Dekker.
- NEWCOMB, S. (1886). A generalized theory of the combination of observations so as to obtain the best result. *Amer. J. Math.* **8**, 343–366.
- PEARSON, K. (1894). Contribution to the mathematical theory of evolution. *Philos. Trans. Roy. Soc. London Ser. A* **185**, 71–110.
- PICKERING, R.M. & FORBES, J.F. (1984). A classification of Scottish infants using latent class analysis. *Statist. Medicine* **3**, 249–259.
- RAO, C.R. (1948). The utilization of multiple measurements in problems of biological classification. *J. Roy. Statist. Soc. Ser. B* **10**, 159–203.
- REDNER, R.A. & WALKER, H.F. (1984). Mixture densities, maximum likelihood and the EM algorithm. *SIAM Rev.* **26**, 195–239.
- SCHNUTE, J. & FOURNIER, D. (1980). A new approach to length–frequency analysis: growth structure. *Canad. J. Fish. Aquat. Sci.* **37**, 1337–1351.
- SCLOVE, S.L. (1977). Population mixture models and clustering algorithms. *Comm. Statist. A — Theory Methods* **6**, 417–434.
- SCOTT, A.J. & SYMONS, M.J. (1971). Clustering methods based on likelihood ratio criteria. *Biometrics* **27**, 387–397.
- SYMONS, M.J. (1981). Clustering criteria and multivariate normal mixtures. *Biometrics* **37**, 35–43.
- TAN, W.Y. & CHAN, W.C. (1972). Some comparisons of the method of moments and the method of maximum likelihood in estimating parameters of a mixture of two normal densities. *J. Amer. Statist. Assoc.* **67**, 702–708.
- TITTERINGTON, D.M., SMITH, A.F.M. & MAKOV, U.E. (1985). *Statistical Analysis of Finite Mixture Distributions*. New York: Wiley.
- VERMUNT, J.K. (1997).  $\ell$ EM. Technical report, Tilburg University.
- WALLACE, C.S. & BOULTON, D.M. (1968). An information measure for classification. *Comput. J.* **11**, 185–194.
- & DOWE, D.L. (1998). MML mixture modelling of multi-state, Poisson, von Mises circular and Gaussian distributions. In *Proc. 28th Symposium on the Interface*, eds L. Billard & N.I. Fisher. *Computing Science and Statistics*, Vol. 28, pp. 608–613. Fairfax Station, VA: Interface Foundation of North America.
- & FREEMAN, P.R. (1987). Estimation and inference by compact coding. *J. Roy. Statist. Soc. Ser. B* **49**, 223–265.
- WERMUTH, N. (1976a). Analogies between multiplicative models in contingency tables and covariance selection. *Biometrics* **32**, 95–108.
- (1976b). Model search among multiplicative models. *Biometrics* **32**, 253–263.
- WOLFE, J.H. (1967). NORMIX: Computations for estimating the parameters of multivariate normal mixtures of distributions. US Naval Personnel Research Activity, Research Memo SRM 68-2, San Diego.
- (1970). Pattern clustering by multivariate mixture analysis. *Multivariate Behavioural Research* **5**, 329–350.
- (1971). A Monte Carlo study of the sampling distribution of the likelihood ratio for mixtures of multinormal distributions. US Naval Personnel and Training Research Laboratory, Technical Bulletin STB 72-2, San Diego.
- ZHAORONG, J., MCGILCHRIST, C.A. & JORGENSEN, M.A. (1992). Mixed model discrete regression. *Biometrical J.* **34**, 691–700.