



Model-Based Gaussian and Non-Gaussian Clustering

Jeffrey D. Banfield, Adrian E. Raftery

Biometrics, Volume 49, Issue 3 (Sep., 1993), 803-821.

Stable URL:

<http://links.jstor.org/sici?sici=0006-341X%28199309%2949%3A3%3C803%3AMGANC%3E2.0.CO%3B2-G>

Your use of the JSTOR archive indicates your acceptance of JSTOR's Terms and Conditions of Use, available at <http://www.jstor.org/about/terms.html>. JSTOR's Terms and Conditions of Use provides, in part, that unless you have obtained prior permission, you may not download an entire issue of a journal or multiple copies of articles, and you may use content in the JSTOR archive only for your personal, non-commercial use.

Each copy of any part of a JSTOR transmission must contain the same copyright notice that appears on the screen or printed page of such transmission.

Biometrics is published by International Biometric Society. Please contact the publisher for further permissions regarding the use of this work. Publisher contact information may be obtained at <http://www.jstor.org/journals/ibs.html>.

Biometrics

©1993 International Biometric Society

JSTOR and the JSTOR logo are trademarks of JSTOR, and are Registered in the U.S. Patent and Trademark Office. For more information on JSTOR contact jstor-info@umich.edu.

©2002 JSTOR

Model-Based Gaussian and Non-Gaussian Clustering

Jeffrey D. Banfield

Department of Mathematical Sciences, Montana State University,
Bozeman, Montana 59715, U.S.A.

and

Adrian E. Raftery

Department of Statistics, GN-22, University of Washington,
Seattle, Washington 98195, U.S.A.

SUMMARY

The classification maximum likelihood approach is sufficiently general to encompass many current clustering algorithms, including those based on the sum of squares criterion and on the criterion of Friedman and Rubin (1967, *Journal of the American Statistical Association* **62**, 1159–1178). However, as currently implemented, it does not allow the specification of which features (orientation, size, and shape) are to be common to all clusters and which may differ between clusters. Also, it is restricted to Gaussian distributions and it does not allow for noise.

We propose ways of overcoming these limitations. A reparameterization of the covariance matrix allows us to specify that some, but not all, features be the same for all clusters. A practical framework for non-Gaussian clustering is outlined, and a means of incorporating noise in the form of a Poisson process is described. An approximate Bayesian method for choosing the number of clusters is given.

The performance of the proposed methods is studied by simulation, with encouraging results. The methods are applied to the analysis of a data set arising in the study of diabetes, and the results seem better than those of previous analyses. A magnetic resonance image (MRI) of the brain is also analyzed, and the methods appear successful in extracting the main features of anatomical interest. The methods described here have been implemented in both Fortran and S-PLUS versions, and the software is freely available through StatLib.

1. Introduction

Cluster analysis has developed mainly through the invention and empirical investigation of ad hoc methods, in isolation from more formal statistical procedures. In recent years it has been found that basing cluster analysis on a probability model can be useful both for understanding when existing methods are likely to be successful, and for suggesting new methods (Symons, 1981; McLachlan, 1982; McLachlan and Basford, 1988).

One such probability model is that the population of interest consists of G different subpopulations, and that the density of a p -dimensional observation \mathbf{x} from the k th subpopulation is $f_k(\mathbf{x}; \boldsymbol{\theta})$ for some unknown vector of parameters $\boldsymbol{\theta}$. Given observations $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$, we let $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_n)^T$ denote the identifying labels, where $\gamma_i = k$ if \mathbf{x}_i comes from the k th subpopulation. In the so-called classification maximum likelihood procedure, $\boldsymbol{\theta}$ and $\boldsymbol{\gamma}$ are chosen so as to maximize the likelihood

$$L(\boldsymbol{\theta}, \boldsymbol{\gamma}) = \prod_{i=1}^n f_{\gamma_i}(\mathbf{x}_i; \boldsymbol{\theta}). \quad (1.1)$$

Scott and Symons (1971) have worked out the solution when $f_k(\mathbf{x}; \boldsymbol{\theta})$ is multivariate normal with mean vector $\boldsymbol{\mu}_k$ and variance matrix $\boldsymbol{\Sigma}_k$, a distribution which we denote by $\text{MVN}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$. When $\boldsymbol{\Sigma}_k = \sigma^2 I$ ($k = 1, \dots, G$), this reduces to the sum of squares criterion (Gordon, 1981, pp. 50–51), whereas when $\boldsymbol{\Sigma}_k = \boldsymbol{\Sigma}$ ($k = 1, \dots, G$) it yields the criterion of Friedman and Rubin (1967). For a more detailed review of these ideas, see Gordon (1981, Chap. 3).

Key words: Bayes factors; Classification; Diabetes; Hierarchical agglomeration; Iterative relocation; Magnetic resonance imaging; Mixture models.

However, as currently implemented, the classification maximum likelihood procedure has several limitations:

1. It considers only the restrictive model where the covariance matrices are constant across all clusters, or the unparsimonious model where they are arbitrary and unequal. The latter is rarely used in practice, probably because of difficulties caused by its very generality and lack of parsimony (Symons, 1981). It would seem desirable to have criteria based on intermediate models that allow some of the characteristics of the covariance matrices to differ across clusters. For example, clusters may be elliptical with roughly the same size and shape, but oriented in different directions.
2. It allows only for Gaussian distributions, whereas other distributions may be more appropriate in some situations. An example of this arises frequently in unsupervised pattern recognition, where edges may be represented by points clustered uniformly, rather than normally, along a straight line (Banfield and Raftery, 1992; Vesecky et al., 1988).
3. It does not, in general, allow for noise, or data points that do not fit the prevailing pattern of clusters. Indeed, when the covariance matrices are unequal, each cluster must contain at least $p + 1$ observations (Symons, 1981).

In this article, we present a framework for model-based clustering that is sufficiently general to overcome these limitations. In Section 2, we develop maximum likelihood criteria for Gaussian clustering that allow clusters to have different orientations or sizes, while preserving some common features, such as shape. In Section 3, we present practical criteria for non-Gaussian clustering, and we extend the framework to incorporate Poisson noise. In Section 4, we present a model-based approximate Bayesian approach to choosing the number of clusters. In Section 5 we report the results of a Monte Carlo study of the methods presented, and in Section 6 we study their performance on two data sets.

2. Allowing Orientation and Size to Vary Between Clusters in the Gaussian Case

When $f_k(\mathbf{x}; \boldsymbol{\theta})$ is a MVN $(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ density, the likelihood (1.1) has the form

$$L(\boldsymbol{\theta}, \boldsymbol{\gamma}) = \text{const.} \prod_{k=1}^G \prod_{i \in E_k} |\boldsymbol{\Sigma}_k|^{-1/2} \exp\{-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}(\mathbf{x}_i - \boldsymbol{\mu}_k)\}, \quad (2.1)$$

where $E_k = \{i: \gamma_i = k\}$. The maximum likelihood estimator of $\boldsymbol{\mu}_k$ is $\bar{\mathbf{x}}_k = n_k^{-1} \sum_{i \in E_k} \mathbf{x}_i$, where n_k is the number of elements in E_k . Replacing $\boldsymbol{\mu}_k$ in (2.1) with the MLE, $\bar{\mathbf{x}}_k$, yields the concentrated log-likelihood

$$l(\boldsymbol{\theta}, \boldsymbol{\gamma}) = \text{const.} - \frac{1}{2} \sum_{k=1}^G \{\text{tr}(W_k \boldsymbol{\Sigma}_k^{-1}) + n_k \log |\boldsymbol{\Sigma}_k| \}, \quad (2.2)$$

where W_k is the sample cross-product matrix for the k th cluster, namely

$$W_k = \sum_{i \in E_k} (\mathbf{x}_i - \bar{\mathbf{x}}_k)(\mathbf{x}_i - \bar{\mathbf{x}}_k)^T.$$

Note that W_k/n_k is the MLE of $\boldsymbol{\Sigma}_k$.

If $\boldsymbol{\Sigma}_k = \sigma^2 \mathbf{I}$ ($k = 1, \dots, G$), then the log-likelihood (2.2) is maximized by choosing $\boldsymbol{\gamma}$ to minimize $\text{tr}(W)$, where $W = \sum_{k=1}^G W_k$. This is the sum of squares criterion which underlies, for example, Ward's (1963) agglomerative hierarchical clustering method. If $\boldsymbol{\Sigma}_k = \boldsymbol{\Sigma}$ ($k = 1, \dots, G$), then the log-likelihood (2.2) is maximized by choosing $\boldsymbol{\gamma}$ to minimize $|W|$, the criterion of Friedman and Rubin (1967). Finally, when the $\boldsymbol{\Sigma}_k$ are not constrained in any way, the likelihood is maximized by choosing $\boldsymbol{\gamma}$ to minimize $\sum_{k=1}^G n_k \log |W_k/n_k|$. This is similar to, but not the same as, equation (14) of Scott and Symons (1971), which we have been unable to reproduce exactly.

Here we develop new criteria that are more general than that of Friedman and Rubin (1967), but based on more parsimonious models than that of Scott and Symons (1971). They allow some but not all of the features of cluster distributions (orientation, size, and shape) to vary between clusters, while constraining others to be the same. The key to this is a reparameterization of the covariance matrix $\boldsymbol{\Sigma}_k$ in terms of its eigenvalue decomposition

$$\boldsymbol{\Sigma}_k = D_k \Lambda_k D_k^T, \quad (2.3)$$

where D_k is the matrix of eigenvectors and Λ_k is a diagonal matrix with the eigenvalues of $\boldsymbol{\Sigma}_k$ on the diagonal. The orientation of the principal components of $\boldsymbol{\Sigma}_k$ is determined by D_k , while Λ_k specifies the size and shape of the density contours. We write $\Lambda_k = \lambda_k A_k$, where λ_k is the first eigenvalue of $\boldsymbol{\Sigma}_k$, $A_k = \text{diag}\{\alpha_{1k}, \dots, \alpha_{pk}\}$, and $1 = \alpha_{1k} \geq \alpha_{2k} \geq \dots \geq \alpha_{pk} > 0$. Thus D_k determines the orientation of

the k th cluster, λ_k its size, and A_k its shape. By size we mean the volume occupied by the cluster in p -space rather than the number of elements it contains. If the α_{jk} 's are of similar magnitude, then the k th cluster will tend to be hyperspherical, whereas if $\alpha_{2k} \ll 1$, it will be concentrated about a line, and if $\alpha_{2k} \approx 1$ and $\alpha_{3k} \ll 1$ it will be concentrated about a two-dimensional plane in p -space, and so on.

This analysis suggests that the sum of squares criterion is likely to be most appropriate when the clusters all have the same dispersion (Symons, 1981). The criterion of Friedman and Rubin (1967), based on the assumption that D_k , λ_k , and A_k are the same, favors clusters that are ellipsoidal with the same orientation and size, whereas the criterion of Scott and Symons (1971), which assumes all the components to be different, favors clusters of different orientations, shapes, and sizes. By allowing some but not all of these quantities to vary between clusters, we obtain criteria that are appropriate for various intermediate situations.

Assuming that $\Sigma_k = \lambda_k I$ leads to a generalization of the sum of squares criterion. The fact that Σ_k is a multiple of the identity matrix indicates that the underlying densities are spherical. Allowing λ_k to vary between densities allows the sizes of the clusters to differ. The resulting criterion to be minimized is

$$\sum_{k=1}^G n_k \log \operatorname{tr} \left(\frac{W_k}{n_k} \right).$$

Our analysis indicates that this criterion will be most appropriate when the clusters are hyperspherical, but of different sizes.

Next, we allow the orientations to vary while keeping size and shape constant, by requiring that $\lambda_k = \lambda$, $A_k = A$ ($k = 1, \dots, G$), where A is known, and by allowing the D_k 's to vary between clusters. By noting that $|\Sigma_k| = \lambda^p \prod_{j=1}^p \alpha_j$, replacing D_k and λ in (2.2) with their maximum likelihood estimators and writing the eigenvalue decomposition of W_k as

$$W_k = L_k \Omega_k L_k^T, \quad (2.4)$$

where $\Omega_k = \operatorname{diag}\{\omega_{1k}, \dots, \omega_{pk}\}$ and ω_{jk} is the j th eigenvalue of W_k , we see that the resulting log-likelihood is maximized by choosing γ to minimize $S = \sum_{k=1}^G S_k$, where $S_k = \operatorname{tr}(A^{-1} \Omega_k)$. When $p = 2$, this is the criterion that underlies the clustering method of Murtagh and Raftery (1984).

We now allow both size, γ_k , and orientation, D_k , to vary between clusters, while assuming that the shape matrix A is constant across clusters. In this setting, the maximum likelihood estimator of γ is obtained by minimizing

$$S^* = \sum_{k=1}^G n_k \log \left(\frac{S_k}{n_k} \right). \quad (2.5)$$

Table 1 shows the relationships among orientation, size, and shape for the criteria we have found to be the most useful in practice. There are, of course, criteria based on other combinations of these factors.

It is usually not feasible to find the global minimum by evaluating the criterion for all possible partitions of the observations. Many algorithms have been devised for finding local minima or good suboptimal solutions, particularly for the sum of squares criterion. These involve agglomeration,

Table 1
Constraints imposed on clusters by different criteria

Criterion	Origin	Distribution	Orientation	Size	Shape
$\operatorname{tr}(W)$	Ward (1963)	Spherical	Undefined	Same	Same
$ W $	Friedman and Rubin (1967)	Ellipsoidal	Same	Same	Same
S	Murtagh and Raftery (1984)	Ellipsoidal	Different	Same	Same
$\sum_{k=1}^G n_k \log \operatorname{tr}(W_k/n_k)$	This article	Spherical	Undefined	Different	Same
S^*	This article	Ellipsoidal	Different	Different	Same
$\sum_{k=1}^G n_k \log \left \frac{W_k}{n_k} \right $	Scott and Symons (1971)	Ellipsoidal	Different	Different	Different

iterative relocation, or other methods; for reviews see Gordon (1981, 1987), Murtagh (1985), and Jain and Dubes (1988). Algorithms developed for the sum of squares criterion can be adapted for use with the other criteria in Table 1. For example, Murtagh and Raftery (1984) showed how Ward's (1963) agglomerative hierarchical method based on the sum of squares criterion can be generalized for use with the criterion S .

3. Non-Gaussian Clustering and Noise

3.1 Non-Gaussian Clustering: The Uniform-Normal Case

The model (1.1) is general enough to encompass clusters with non-Gaussian distributions. To date, attention has been focused on the multivariate normal distribution because it leads to relatively simple criteria. Here we suggest practical criteria for some non-Gaussian situations.

The basic idea is the use of a local parameterization. We assume that there are matrices D_k ($k = 1, \dots, G$) such that if $\mathbf{z}_i = D_{\gamma_i}(\mathbf{x}_i - \boldsymbol{\mu}_{\gamma_i})$, then \mathbf{z}_i has the density $g_{\gamma_i}(\mathbf{z}_i; \boldsymbol{\theta})$; often these densities will be the same, perhaps modulo a scale parameter. In this general framework, criteria can be derived by maximizing the likelihood, as in Section 2. When the distribution of \mathbf{x}_i is $\text{MVN}(\boldsymbol{\mu}_{\gamma_i}, \boldsymbol{\Sigma}_{\gamma_i})$, and D_k is defined by (2.3), then \mathbf{z}_i is the value of \mathbf{x}_i in the local coordinate system with origin at $\boldsymbol{\mu}_{\gamma_i}$ and axes along the principal components of $\boldsymbol{\Sigma}_{\gamma_i}$.

We now carry out a more detailed analysis of one specific, but important, non-Gaussian situation—when observations are clustered uniformly along and tightly about a line segment in p -space. Such situations arise in ecology when the data include the geographic locations of plants or animals that may be clustered about roughly linear natural features such as rivers or valleys. They also arise in unsupervised pattern recognition, where observations may be edge elements in an image, or data points in a point pattern with a linear feature.

We let $u_i = z_{i1}$, and $\mathbf{v}_i = (v_{i1}, \dots, v_{i,p-1})^T = (z_{i2}, \dots, z_{ip})^T$. We assume that u_i is uniformly distributed between $\phi_{\gamma_i,1}$ and $\phi_{\gamma_i,2}$, and that $\mathbf{v}_i \sim \text{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_k)$. Let $\phi_k = \phi_{k2} - \phi_{k1}$ and $\boldsymbol{\Sigma}_k = \sigma_k^2 I$; typically σ_k will be small compared to ϕ_k .

We now derive an approximate maximum likelihood estimator for $\boldsymbol{\gamma}$ under this model. We estimate D_k by $\hat{D}_k = L_k$, where L_k is defined by (2.4), and we estimate $\boldsymbol{\mu}_k$ by $\hat{\boldsymbol{\mu}}_k = \bar{\mathbf{x}}_k$. We then define $\hat{\mathbf{z}}_i = \hat{D}_{\gamma_i}(\mathbf{x}_i - \bar{\mathbf{x}}_{\gamma_i})$, with corresponding definitions for \hat{u}_i and $\hat{\mathbf{v}}_i$. Conditionally on these estimated values of D_k and $\boldsymbol{\mu}_k$, the log-likelihood for $\boldsymbol{\phi} = (\phi_1, \dots, \phi_G)^T$, σ_k^2 , and $\boldsymbol{\gamma}$ is

$$l(\boldsymbol{\phi}, \sigma^2, \boldsymbol{\gamma}) = \text{const.} - \sum_{k=1}^G \left\{ n_k \log \phi_k + \frac{1}{2} (p-1) n_k \log \sigma_k^2 + \frac{1}{2\sigma_k^2} \sum_{i \in E_k} \hat{\mathbf{v}}_i^T \hat{\mathbf{v}}_i \right\}. \quad (3.1)$$

If we assume that $\sigma_k^2 = \sigma^2$ ($k = 1, \dots, G$), then the log-likelihood in equation (3.1) is maximized by

$$\hat{\phi}_k = \max_{i \in E_k} \{\hat{u}_i\} - \min_{i \in E_k} \{\hat{u}_i\}, \quad \hat{\sigma}^2 = \{n(p-1)\}^{-1} \sum_{i=1}^n \hat{\mathbf{v}}_i^T \hat{\mathbf{v}}_i.$$

We therefore choose $\boldsymbol{\gamma}$ to minimize the criterion

$$\frac{1}{2} (p-1)n \log \hat{\sigma}^2 + \sum_{k=1}^G n_k \log \hat{\phi}_k. \quad (3.2)$$

In the situation where the σ_k^2 's are not constant across clusters we obtain

$$\hat{\sigma}_k^2 = \{n_k(p-1)\}^{-1} \sum_{i \in E_k} \hat{\mathbf{v}}_i^T \hat{\mathbf{v}}_i,$$

and $\boldsymbol{\gamma}$ is chosen to minimize the criterion

$$U = \sum_{k=1}^G \left\{ \frac{1}{2} (p-1)n_k \log \hat{\sigma}_k^2 + n_k \log \hat{\phi}_k \right\}. \quad (3.3)$$

Many variations on this "uniform-normal" theme are possible, and lead to simple criteria. For example, clusters may be distributed tightly about a two-dimensional planar region in p -space; this can be represented by specifying the distribution of (z_{i1}, z_{i2}) to be concentrated on such a region. Also, the distribution of the scatter about the main line segment or planar region may be more general than assumed above, leading, for example, to a range of values for the covariance matrix of \mathbf{v}_i , such as those considered in Section 2.

3.2 Allowing for Noise

So far, we have assumed that each observation belongs to a cluster. However, even when a data set is made up mainly of clusters of the prescribed type, there may be other data points that do not follow

this pattern. We allow for this possibility by extending our model to include such observations, assumed to arise from a Poisson process with intensity ν . Let $E = \bigcup_{k=1}^G E_k$ and $n_0 = n - \sum_{k=1}^G n_k$. Then the likelihood (1.1) is modified as follows:

$$L(\boldsymbol{\theta}, \nu, \boldsymbol{\gamma}) = \frac{(\nu A)^{n_0} e^{-\nu A}}{n_0!} \prod_{i \in E} f_{\gamma_i}(\mathbf{x}_i; \boldsymbol{\theta}), \quad (3.4)$$

where A is the hypervolume of the region from which the data have been drawn. The clustering criteria developed so far can easily be modified to be based on (3.4). Taking account of noise in this way facilitates our proposed method for choosing the number of clusters, described in Section 4.

4. Choosing the Number of Clusters: An Approximate Bayesian Approach

Here we suggest an approximate Bayesian approach to the choice of the number of clusters. We first write down an exact Bayesian solution, but this usually cannot be computed in a reasonable amount of time. Arguing heuristically, we obtain an approximation to the Bayesian solution that seems to work well in numerical experiments, some of which are reported in Sections 5 and 6.

We view the problem of estimating the number of clusters as one of choosing between competing models for the same data. The exact Bayesian solution consists of finding the posterior probability $p(G | \mathbf{x})$ of each number of clusters G given the data \mathbf{x} . This approach seems to have advantages over the alternative of hypothesis testing in the general context of model comparison, as it avoids the problems of multiple comparisons, comparing nonnested models, and the tendency of hypothesis tests to select unparsimonious models when the sample size is large (Berger and Sellke, 1987; Raftery, 1986b, and Technical Report No. 121, Department of Statistics, University of Washington, 1988). The details have been worked out for many statistical problems, including the general linear model (Smith and Spiegelhalter, 1980), generalized linear models (Raftery, 1986a, and unpublished technical report cited above), change-points and point processes (Akman and Raftery, 1986; Raftery and Akman, 1986), and software reliability (Raftery, 1987, 1988).

Technically, it is easiest to start with the Bayes factor, or ratio of posterior to prior odds for $G = r$ against $G = s$, defined by

$$B_{rs} = p(\mathbf{x} | G = r) / p(\mathbf{x} | G = s). \quad (4.1)$$

In (4.1), $p(\mathbf{x} | G = r)$ is the marginal likelihood

$$p(\mathbf{x} | G = r) = \sum_{\boldsymbol{\gamma} \in \Gamma_r} \int \int L(\boldsymbol{\theta}, \nu, \boldsymbol{\gamma}) p(\boldsymbol{\theta}, \nu, \boldsymbol{\gamma}) d\boldsymbol{\theta} d\nu,$$

where $\Gamma_r = \{0, 1, \dots, r\}^n$, $L(\boldsymbol{\theta}, \nu, \boldsymbol{\gamma})$ is the generalized likelihood defined by (3.4), and $p(\boldsymbol{\theta}, \nu, \boldsymbol{\gamma})$ is the joint prior density of $\boldsymbol{\theta}$, ν , and $\boldsymbol{\gamma}$. When $\gamma_i = 0$, \mathbf{x}_i belongs to the “noise” and appears in the Poisson part of the likelihood (3.4). A different but related approach is described in unpublished work by Rissanen.

Here we concentrate on the approximate calculation of $B_{r,r+1}$ ($r = 1, \dots, n - 1$). This yields posterior probabilities $p(G = r | \mathbf{x})$ directly, as follows. Noting that $B_{rs} = \prod_{t=1}^{s-r} B_{r+t-1, r+t}$ ($r < s$) and $B_{sr} = B_{rs}^{-1}$, we calculate B_{rs_0} for $r = 1, \dots, n - 1$ and some fixed s_0 . Then

$$p(G = r | \mathbf{x}) = B_{rs_0} p(G = r) / \sum_{t=1}^{n-1} B_{ts_0} p(G = t), \quad (4.2)$$

where $p(G = r)$ is the prior probability that there are r clusters.

We approximate $B_{r,r+1}$ as follows. In an agglomerative hierarchical clustering algorithm, the choice between $G = r + 1$ and $G = r$ is a decision whether or not to merge two particular clusters into one. In the p -dimensional multivariate normal case, this is exactly a standard comparison of nested hypotheses in the general linear model, and Smith and Spiegelhalter (1980) have shown that in that case minus twice the logarithm of the Bayes factor is approximately equal to

$$\lambda_r - \left\{ \frac{3}{2} + \log(p n_{r,r+1}) \right\} \delta_r, \quad (4.3)$$

where λ_r is the likelihood ratio test statistic, δ_r is the number of degrees of freedom in the asymptotic chi-square distribution of λ_r , and $n_{r,r+1}$ is the number of observations in the merged cluster. However, (4.3) is invalid in the clustering context because the regularity conditions on which it is based do not hold. Wolfe (1971) suggested getting around the problem by doubling the number of degrees of freedom, and Hernandez-Alvi (unpublished Ph.D. thesis, University of Oxford, 1979) found that to be a reasonable approximation. Aitkin, Anderson, and Hinde (1981) and McLachlan (1987) had some reservations about the use of Wolfe’s (1971) approximation when δ_r is large, but the simulations

of Everitt (1981) showed it to perform well for values of δ_r between 1 and 5, which is the range of primary interest to us. We therefore use the approximation

$$-2 \log B_{r,r+1} \approx \lambda_r - \left\{\frac{3}{2} + \log(pn_{r,r+1})\right\}2\delta_r = E_r, \tag{4.4}$$

where δ_r is now the decrease in the number of parameters caused by going from $G = r + 1$ to $G = r$. In Table 2, for the case where the data are two-dimensional, we show the values of δ_r and the individual cluster parameters that must be estimated for the clustering criteria from Sections 2 and 3. We write

$$D = \begin{pmatrix} \cos \Psi & -\sin \Psi \\ \sin \Psi & \cos \Psi \end{pmatrix},$$

where Ψ is the orientation of the cluster. For the criteria in Section 3, ϕ_k can be superefficiently estimated, i.e., the asymptotic variance goes to zero faster than the usual rate of $O(1/n)$ [here it is $O(1/n^2)$], and so it is not included in the count. For the models corresponding to all the criteria in Table 1 except that of Friedman and Rubin (1967), the term λ_r in (4.4) involves only the contributions to the likelihood of the clusters involved in the merger. If we define the maximized log likelihoods for the two clusters that are merged as $l_{k'}$ and $l_{k''}$ and the maximized likelihood for the cluster resulting from the merger as l_k , we may write

$$\lambda_r = 2(l_{k'} + l_{k''} - l_k). \tag{4.5}$$

The likelihoods for the clusters that are not involved in the merger cancel out in the likelihood ratio.

If we assume that the clusters are embedded in a Poisson process, the outcomes of the mergers are slightly more complicated since at each stage in the agglomerative process the number of clusters, G , can increase, decrease, or remain the same. The reason for this is that we have two types of data: clusters and noise. If we form a new cluster by merging two singletons that were considered noise, then G will increase. If we merge a singleton with an existing cluster, then G will not change. If two existing clusters are merged, then G will decrease. If two singletons are merged to form cluster k , then $\lambda_r = 2l_k$, and δ_r for the merger is equal to minus the value of δ_r given in Table 2. If a singleton is merged with cluster k' to form cluster k then $\lambda_r = -2(l_{k'} - l_k)$ and $\delta_r = 0$ since the parameterization has not changed. When two existing clusters, k' and k'' , are merged to form cluster k , λ_r is given by equation (4.5) and δ_r is as given in Table 2.

Having obtained $B_{r,r+1}$ ($r = 1, \dots, n - 1$) from (4.4), we may calculate $p(G = r | \mathbf{x})$ ($r = 1, \dots, n - 1$) using (4.2). A simple approach is to use as an estimate of the number of clusters the value of r for which $p(G = r | \mathbf{x})$ is greatest. However, (4.4) provides a rather crude approximation to $p(G = r | \mathbf{x})$. We therefore prefer to consider several values of the number of clusters, guided by the values of $p(G = r | \mathbf{x})$, or, equivalently, by F_r , defined by $F_1 = 0$ and $F_r = \sum_{i=1}^{r-1} E_i \approx \text{constant} + 2 \log p(G = r | \mathbf{x})$ ($r \geq 2$). Following Good (1983), we refer to F_r as the *approximate weight of evidence* (AWE) for the number of clusters being r . In our experience, the change in the approximate weight of evidence, $E_r = F_{r+1} - F_r$, is often large and positive for the first few values of r ($r = 1, \dots, R$, say) and small or negative thereafter. If that is the case, considerations of parsimony lead us to consider $G = R + 1$, as well as the value of r that maximizes the approximate weight of evidence, and any intervening values.

Table 2
Decrease in the number of parameters caused by reducing the number of clusters by one, for several criteria, in the two-dimensional case

Criterion	δ_r	Parameters
$\text{tr}(W)$	2	μ_x, μ_y
$ W $	2	μ_x, μ_y
S	3	μ_x, μ_y, Ψ
$\sum_{k=1}^G n_k \log \text{tr}(W_k/n_k)$	3	μ_x, μ_y, λ_k
S^*	4	$\mu_x, \mu_y, \Psi, \lambda_k$
$\sum_{k=1}^G n_k \log \left \frac{W_k}{n_k} \right $	5	$\mu_x, \mu_y, \Psi, \lambda_k, \alpha_{2k}$
Equation (3.2)	3	μ_x, μ_y, Ψ
U	4	$\mu_x, \mu_y, \Psi, \sigma_k^2$

5. Simulations

5.1 Simulated Clusters

Figure 1(a) shows three clusters generated from bivariate normal distributions with the same shape but different sizes and orientations. Figures 1(b), 1(c), and 1(d) show the results of grouping the data into three clusters using the criteria S^* , $\text{tr}(W)$, and single-link method (SL), respectively. The S^* criterion performed well. Three of the four misclassified points are within or close to the intersections of the clusters. This is inevitable, since even the human eye, with its remarkable pattern recognition and classification abilities, finds it hard to classify points at the intersection of clusters. Ward's criterion, $\text{tr}(W)$, misclassified 18 of the 45 points and did not reproduce the general shape of the clusters. As can be seen from Figure 1(c), it tends, instead, to find "circular" clusters. The single-link method has been suggested for finding long clusters such as those in Figure 1(a). However, as can be seen from Figure 1(d) and Tables 3, 4, and 5, it does not perform well when the clusters intersect.

Clusters that are physically separate, in whatever metric is being used, are easy to distinguish with most clustering criteria. The clusters we have been working with are distinguished from each other by their structure. A point within one cluster may be closer, in Euclidean distance, to points in other clusters than to any point in the cluster to which it belongs, yet we are able to classify it correctly due to the known structure of the clusters. For example, consider Figure 1(b). Note the two points on the left that have been correctly classified as belonging to cluster 2 (triangles), yet they are closer to points

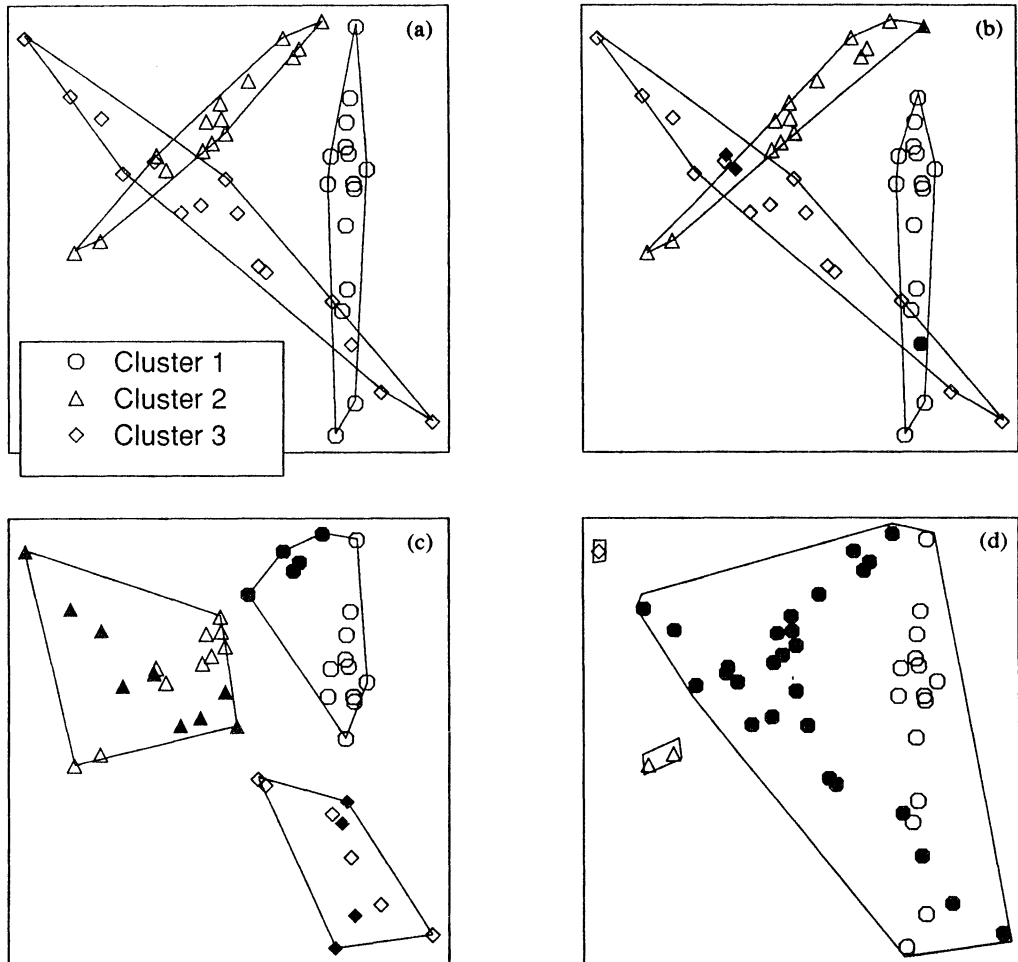


Figure 1. (a) Three clusters generated from bivariate normal distributions with the same shape but different sizes and orientations. The solid lines are the convex hulls of the groups. (b) The clusters formed by the S^* criterion. The filled-in symbols represent misclassified points. For example, the filled-in triangle at the top right-hand corner was classified as a diamond, but in fact is a circle. (c) The clusters found by Ward's sum-of-squares criterion, $\text{tr}(W)$. (d) The clusters found by the single-link method.

in cluster 3 (diamonds) than to any point in cluster 2. Criteria based strictly on distance measures, such as the single-link method, are unable to handle clusters that are defined by their structure.

5.2 Simulated Clusters with Noise

Figure 2 shows three clusters with added noise. The clusters were generated from bivariate normal distributions with the same shape but different sizes and orientations while the noise was generated by a Poisson process. This example differs from that in Section 5.1 in that noise has been added and we do not assume the numbers of clusters to be known in advance.

After clustering the data in Figure 2 using S^* in a hierarchical agglomeration procedure, the approximate weight of evidence was calculated at each iteration, as shown in Figure 3. The AWE is maximized at seven clusters and falls off sharply after that, indicating that the clustering algorithm should be stopped at seven clusters. Figure 4 shows the results at seven clusters after using an iterative

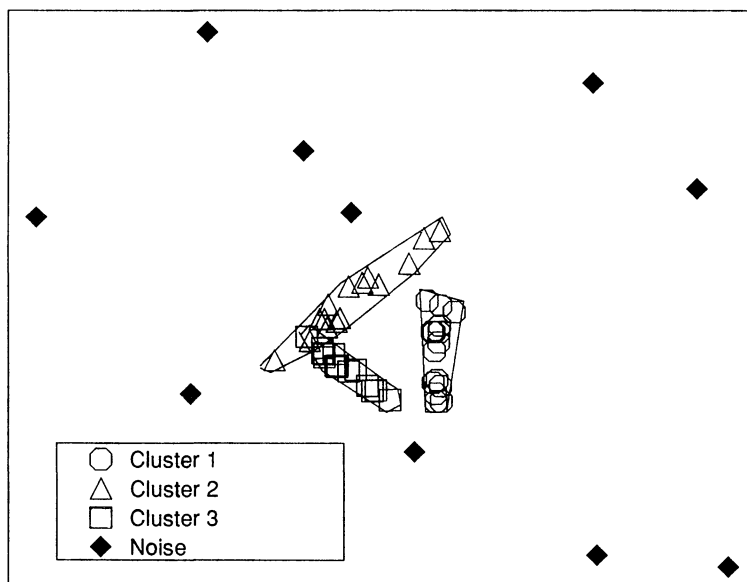


Figure 2. Three clusters with noise. The clusters were generated from bivariate normal distributions with the same shape but different sizes and orientations. The noise was generated from a Poisson process.

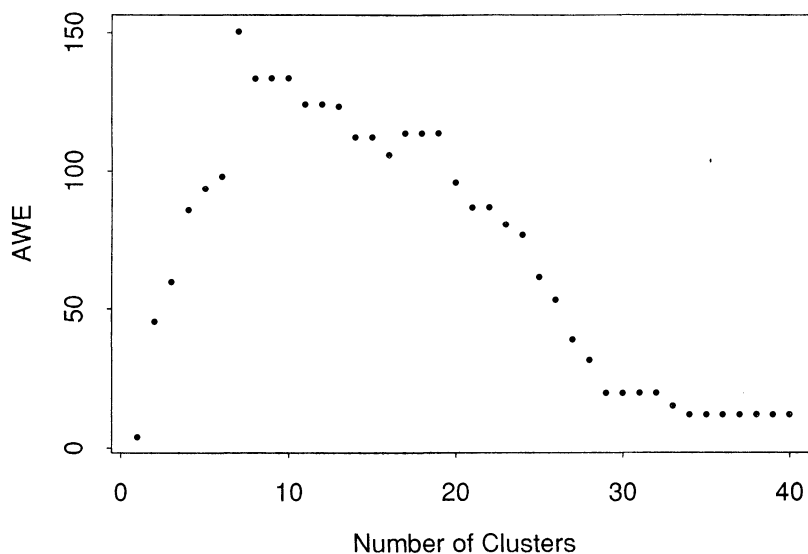


Figure 3. Approximate weight of evidence (AWE) for the number of clusters in Figure 2 using the criterion S^* . The maximum occurs at seven clusters and leads to the results shown in Figure 4.

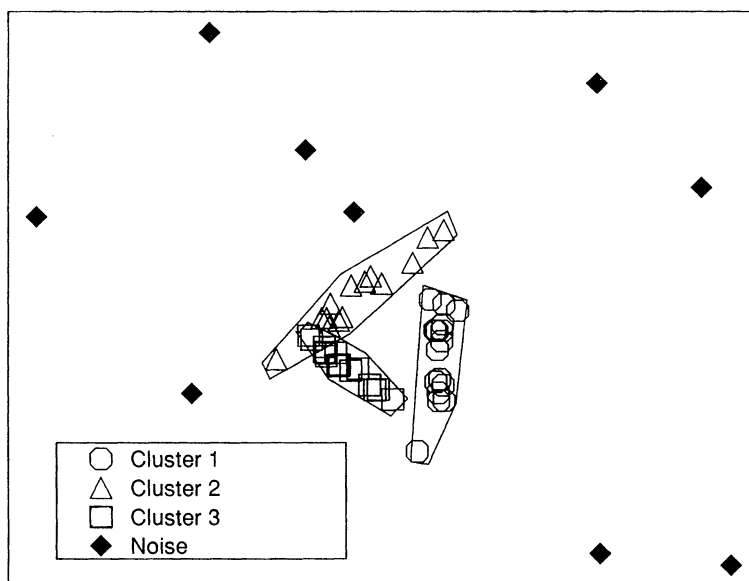


Figure 4. The clusters resulting from the data in Figure 2 using the criterion S^* and stopping with seven clusters as indicated by the \diamond points in Figure 3. Four of the seven clusters had three or fewer points and are considered noise.

relocation algorithm to improve upon the original agglomerative results. Four of the seven clusters contained a total of ten points and have been classified as noise. The three main clusters are well defined with one misclassification, and only one of the noise points has been misclassified.

5.3 Simulation Results

To compare the performance of our clustering criteria with that of standard, commonly used criteria, we carried out a Monte Carlo study. The standard criteria used were the single-link method (SL) and Ward's sum of squares criterion, $\text{tr}(W)$. These were compared with the three criteria S , introduced for two dimensions by Murtagh and Raftery (1984) and generalized to higher dimensions in Section 2, S^* defined by (2.5), and U defined by (3.3).

To compare the criteria we generated 100 random samples from each of the three types of data for each of four values of $\alpha = \lambda_2/\lambda_1$, where the λ_i are the eigenvalues of the covariance matrix, giving a total of 1,200 samples. The three types of data correspond to the three models for which S , S^* , and U are optimal criteria. When generating the data for which U was optimal, ϕ_k was generated from a $U(.2, .6)$ distribution and $\hat{\sigma}_k^2$ was proportional to $\alpha\phi_k^2$. Each sample consisted of three clusters, the orientation of each cluster was randomly chosen from a $U(0, \pi)$ distribution, and the centers were randomly chosen in the unit square. The number of points in each cluster was generated from a discrete uniform distribution on the integers between 15 and 25.

Tables 3, 4, and 5 show the proportion of points misclassified by each of the five criteria considered. The single-link method performed poorly, while Ward's sum of squares did only slightly better. The three criteria S , S^* , and U all performed much better. Of these three, S^* did marginally better than the others, but the differences between them were small. As one might expect, each of the three

Table 3

Bivariate normal clusters with the same size and shape but different orientations. S is the optimal criterion. One hundred random samples were generated for each value of α . The entries in the table are the percentages of points misclassified.

Criterion	α			
	.001	.005	.01	.025
S	4	13	16	25
S^*	4	16	18	24
U	7	19	26	30
SL	51	51	52	52
$\text{tr}(W)$	40	41	41	40

Table 4

Bivariate normal clusters with the same shape but different sizes and orientations. S^ is the optimal criterion. One hundred random samples were generated for each value of α . The entries in the table are the percentages of points misclassified.*

Criterion	α			
	.001	.005	.01	.025
S	14	17	24	31
S^*	7	12	17	22
U	7	12	18	26
SL	46	47	50	48
$\text{tr}(W)$	48	48	46	46

Table 5

Bivariate uniform-normal clusters. The observations are clustered uniformly along and tightly about a line segment in two-dimensional space, as described in Section 3.1. U is the optimal criterion. One hundred random samples were generated for each value of α . The entries in the table are the percentages of points misclassified.

Criterion	α			
	.001	.005	.01	.025
S	4	11	13	18
S^*	5	9	12	19
U	3	7	9	14
SL	38	41	45	43
$\text{tr}(W)$	43	41	44	43

criteria S , S^* , and U performed best on the type of data for which it was designed, but it also performed well on the other kinds of data.

The clear superiority of S , S^* , and U to the single-link and Ward's method held for each combination of the three kinds of data with the four values of α . The results for the three kinds of data were quite similar. As α increased, the proportion of points misclassified by S , S^* , and U increased. This reflects the fact that as α increases, the data-generating mechanism more closely approximates that for which $\text{tr}(W)$ is the best criterion, and so the superiority of S , S^* , and U becomes less marked. Averaged over the 1,200 random samples generated, the proportion of points misclassified was 16% for S , 14% for S^* , 15% for U , 47% for the single-link method, and 43% for Ward's sum of squares.

It is assumed that some prior information about α is available. This can come from a training sample or knowledge of the mechanism generating the data—for example, the resolution of the edge detector used in processing a digital image. Our numerical work, including the examples in Section 6, indicates that our criteria are not overly sensitive to errors in the estimation of α . In the simulation study the correct value of α was used in S , S^* , and U . This provides information on the best performance that can be expected.

6. Examples

6.1 Example 1: Diabetes Data

Reaven and Miller (1979) described and analyzed data consisting of the area under a plasma glucose curve (glucose area), the area under a plasma insulin curve (insulin area) and steady-state plasma glucose response (SSPG) for 145 subjects. The subjects were clinically classified into three groups, chemical diabetes (Type 1), overt diabetes (Type 2), and normal (nondiabetic). Symons (1981) reanalyzed the data using seven different clustering criteria. Taking the clinical classification to be correct, we evaluate one of our criteria and compare it with those considered by Symons (1981), using the data as published in Andrews and Herzberg (1985).

Reaven and Miller (1979, Figs 1–4) showed four two-dimensional projections of the data. The data have the three-dimensional shape of a boomerang with two wings and a fat middle. One of the wings corresponds to patients with overt diabetes, the other wing is composed primarily of patients with chemical diabetes, and the “fat middle” is composed of normal patients. By viewing the data using a rotating three-dimensional scatterplot, such as the ones provided in MacSpin (Donoho, Donoho, and Gasko, 1988) or XLISP-STAT (Tierney, Technical Report No. 528, School of Statistics, University

of Minnesota, 1988), this structure is obvious and several other features become apparent. One of the “wings” is almost planar, the other is linear with some curvature, and the “fat middle” has a shape similar to an American football. Four two-dimensional projections of the data are shown in Figure 5.

Based on this information, we could use the approach developed in Sections 2 and 3 to design a purpose-built clustering criterion for this application. However, we prefer to use a very general criterion of the form S^* , where $A_k = \text{diag}\{1, \alpha, \alpha\}$. This criterion is optimal for trivariate normal clusters with different sizes and orientations but the same “tubular” shape, clustered circularly about a line in \mathbb{R}^3 . The estimated values of α for the three clinically identified groups are .09, .19, and .34. The results were relatively insensitive to changes in α so long as it remained in that broad range. The results we report are for $\alpha = .2$. We have scaled each variable by dividing it by its maximum value.

Starting from the correct clinical classification and using a single point iterative relocation algorithm with the criterion S^* , the optimal classification, as given in Table 6, resulted in only 10% of the points being misclassified. This compares favorably with the results given by all seven clustering criteria used by Symons (1981) for this data set. We also used a hierarchical agglomerative clustering algorithm followed by iterative relocation. Once again, the results compare favorably with those of Symons (1981). The clusters found are shown in Figure 6.

The AWE for the hierarchical agglomerative clustering algorithm increased steadily until the final four iterations. Figure 7 shows the number of clusters versus the AWE over the last 20 iterations. From this it can be seen that the AWE increases sharply as one goes from one cluster to two, and again from two to three. It increases slightly as the number of clusters goes up to four and decreases

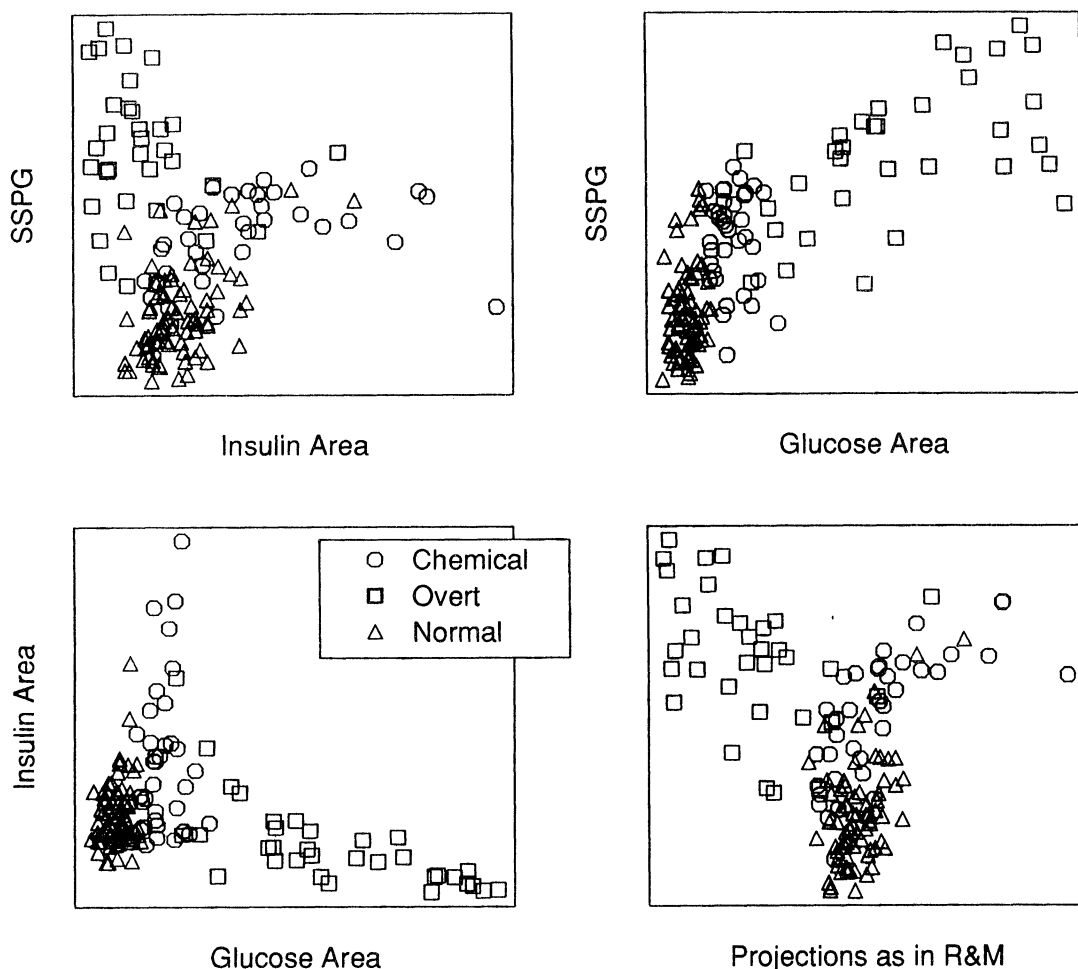


Figure 5. Four two-dimensional projections of the three-dimensional diabetes data of Reaven and Miller (1979). The symbols indicate the clinical classification of subjects as having chemical diabetes, overt diabetes, or being normal. Lower right panel shows the approximate projections represented by the artist’s sketch in Reaven and Miller (1979) and reproduced in Symons (1981).

Table 6

Results of clustering the diabetes data. The first row shows the result of single point iterative relocation using the criterion S^* with $A_k = \text{diag}\{1, .2, .2\}$, starting with the clinical classification. The second row shows the result of hierarchical agglomeration followed by iterative relocation with the same criterion. The remaining rows show the results of seven other clustering procedures, starting at the clinical classification, as reported by Symons (1981). Criterion (13) of Symons (1981) is due to Maronna and Jacovkis (1974). The error rate % is the percentage of the subjects who were not classified in the same way by the clustering method as by the clinical diagnosis.

Method	Error rate %	Clinical classification		
		Normal (76, 0, 0)	Chemical (0, 36, 0)	Overt (0, 0, 33)
S^* from clinical	10	(65, 0, 0)	(11, 36, 4)	(0, 0, 29)
S^* agglomerative	10	(65, 0, 0)	(11, 36, 4)	(0, 0, 29)
$ W $	19	(73, 17, 3)	(3, 19, 4)	(0, 0, 26)
Reaven and Miller (1979) variant of $ W $	14	(73, 10, 1)	(3, 26, 6)	(0, 0, 26)
(8) in Symons (1981)	26	(75, 30, 6)	(1, 6, 1)	(0, 0, 26)
(10) in Symons (1981)	26	(75, 30, 6)	(1, 6, 1)	(0, 0, 26)
(13) in Symons (1981)	13	(73, 10, 0)	(3, 26, 7)	(0, 0, 26)
(11) in Symons (1981)	14	(63, 0, 0)	(13, 30, 2)	(0, 6, 31)
(12) in Symons (1981)	13	(73, 9, 0)	(3, 27, 7)	(0, 0, 26)

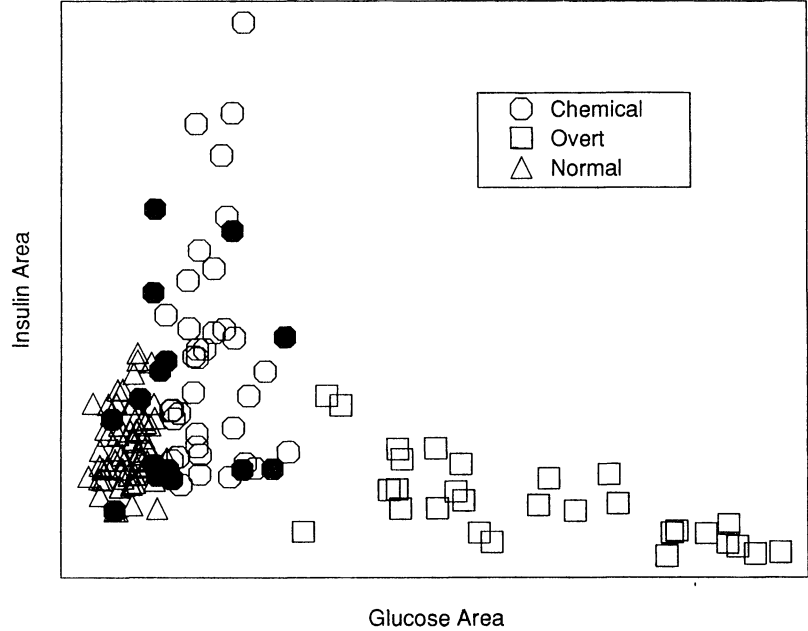


Figure 6. The three clusters in the diabetes data found by hierarchical agglomeration followed by iterative relocation using the criterion S^* with $A_k = \text{diag}\{1, .2, .2\}$. The two-dimensional projection shown is that of Figure 5(c). The symbols indicate the classification of the subjects based on the clustering algorithm. The filled-in symbols represent subjects whose clustering classification differs from the clinical classification.

thereafter, although very slowly until six clusters. If we did not know the true number of clusters this would lead us to focus attention on the groupings into three, and four clusters, and to perform a more detailed analysis on these sets of clusters, keeping in mind that a solution with five or six clusters may be reasonable.

6.2 Example 2: MRI Brain Scan

Magnetic resonance imaging (MRI) is a method for measuring the chemical characteristics of body tissue based on the magnetic resonance of hydrogen nuclei within the tissue (Oldendorf and Oldendorf,

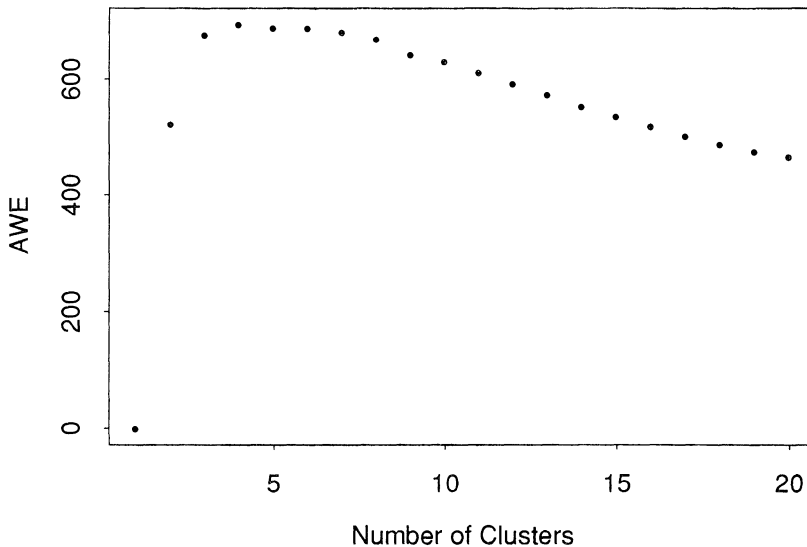


Figure 7. Approximate weight of evidence (AWE) for the number of clusters in the diabetes data over the last 20 iterations of the clustering algorithm. The AWE increases sharply up to three clusters, with a further slight increase going to four clusters. This would lead us to focus on the groupings into three and four clusters.

1988). MRI can unobtrusively distinguish between different tissue types within small-volume elements, called voxels, over very thin cross-sections of the body. Depending on how the hydrogen nuclei are stimulated, different characteristics within each voxel can be measured. This leads to multiple measurements on each voxel within the image, analogous to multiband satellite images such as LANDSAT (an abbreviation for Land Satellite, a series of satellites designed to monitor the Earth's surface). Figure 8 shows three bands from an MRI scan of a human brain.

Problems in MRI scanning include the identification of anatomical structures, the identification of different anatomical structures with similar chemical characteristics (especially when dealing with new organisms), volumetric measurement of specific components (such as fat or gray matter), and smooth rendering of anatomical features for use in creating three-dimensional images of the brain. In this example we provide an initial step in addressing all of these questions by using cluster analysis to identify individual anatomical structures in the phase space. We define the phase space to be the p -dimensional space defined by the different MRI bands without reference to the spatial information. In this example the phase space is three-dimensional.

Since the different MRI bands measure chemical characteristics within each voxel, it is reasonable to assume that voxels with similar chemical components will cluster together in the phase space. Furthermore, there are physical reasons to expect the clusters to have linear shapes since the response

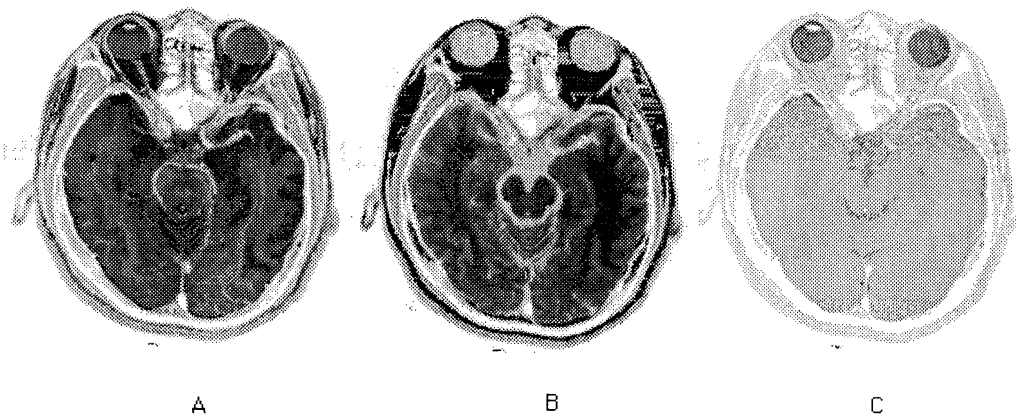


Figure 8. Three bands from a magnetic resonance image (MRI) of the brain. A is the T_1 band, B is the T_2 band, and C is the PD band. Each band contains a circular image but the low-intensity pixels, due to air around the skull, blend into the background (see Figures 11 and 12).

within each voxel can be modeled as a linear combination of the individual chemical components that make up the voxel. This linear mixing model has been used with success in other imaging applications (Adams, Smith, and Johnson, 1986; McDonald and Willis, unpublished videotape, Department of Statistics, University of Washington, 1987).

It is not practical to try to cluster all 26,100 voxels shown in Figure 8. We therefore propose to cluster a random sample of voxels from the image and use the resulting clusters to classify the remaining image voxels. This procedure has the flavor of discriminant analysis except that in discriminant analysis the features of interest must be known in advance. Our approach requires no such specific prior knowledge.

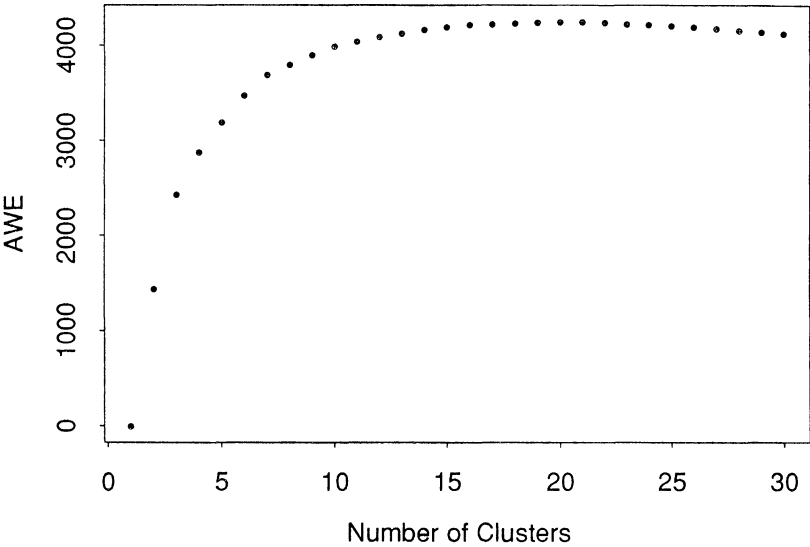


Figure 9. The approximate weight of evidence (AWE) for the MRI brain scan data. The AWE is maximized at 21 clusters.

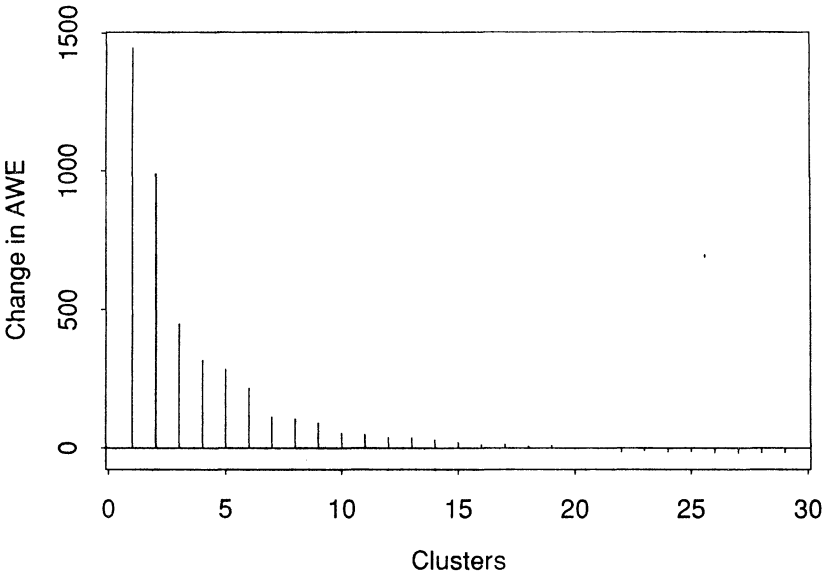


Figure 10. The rate of change in the AWE from Figure 9. For example, the bar above 5 indicates the change in the AWE when going from five to six clusters. It is clear that the AWE increases rapidly from one cluster up to seven clusters, where the rate of increase slows down dramatically. Beyond 20 clusters the changes in the AWE are negative, leading us to explore in more detail the solutions having from 7 to 20 clusters.

We clustered 522 randomly chosen voxels (2% of the total) from the images in Figure 8 using S^* with $A_k = \text{diag}(1, .3, .3)$. We find S^* , with A_k of this form, to be adaptable and robust to the choice of the values for A_k .

We consider numbers of clusters between the number at which the AWE first levels off and the

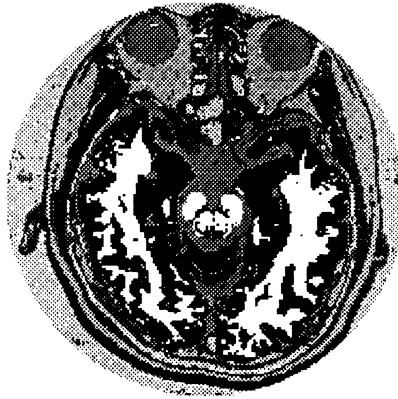


Figure 11. The spatial configuration of the clusters in the seven-cluster solution (the smallest number of clusters indicated by the AWE, Figures 9 and 10). The light gray pixels around the outside of the skull (giving a circular appearance) represent air.



Figure 12. The clusters from Figure 11 shown individually. A: Bone; B: Air; C: White matter; D: Fluids; E: Muscle; F: Fat; G: Gray matter.

number of clusters for which the AWE is maximized. For distinct clusters the AWE indicates an obvious optimal number of clusters as shown in the earlier simulations. When the clusters are not distinct, as is often the case, the AWE provides a reasonable subset. Our experience is that out of this subset the solutions with fewer clusters generally provide more interpretable results.

Figure 9 shows the AWE for the 522 points we have clustered for this example. The maximum occurs at 21 clusters. Figure 10 shows the rate of change in the AWE and indicates that for fewer than seven clusters the fit is seriously degraded. This leads us to consider the solutions that have between 7 and 21 clusters. These results agree with those of independent medical analysis. Six or fewer clusters lead to the merging of anatomically dissimilar structures. Between 7 and 21 clusters most of the mergers combine obviously similar anatomical structures except in the situations where it is uncertain whether MRI can distinguish between structures such as subcutaneous fat and bone marrow or spinal and corneal fluid.

Figure 11 shows the spatial configuration of the seven-cluster solution. Figure 12 shows the individual clusters in Figure 11, together with the anatomical features to which they correspond. These figures, together with the independent medical assessment, suggest that the linear clusters found in the phase space do correspond to distinct anatomical structures. We have used other clustering methods on this problem (for example, complete link clustering and Ward's method) and find that they tend to group together dissimilar anatomical voxels. The reason for this is that the features in the phase space are defined by their linear shape and not simply by Euclidean distance so that criteria based strictly on a distance measure will not be able to identify individual features accurately.

7. Discussion

We have proposed ways of overcoming some of the limitations of the classification maximum likelihood procedure for cluster analysis, as currently implemented. These are (1) the inability to specify some but not all features (orientation, size, shape) to be constant across clusters; (2) the restriction to normal distributions; and (3) the failure to account for "noise." We have also proposed an approximate Bayesian solution to the problem of choosing the number of clusters, which seems to avoid some of the difficulties associated with solutions to this problem based on significance testing.

In the context of Gaussian clustering, we reparameterize the covariance matrices in terms of their eigenvalue decompositions. Each group of parameters then corresponds clearly to a particular feature of the cluster (orientation, size, or shape), and criteria appropriate for a range of different situations result by constraining none, some, or all of these features to be constant across clusters. This leads to a range of criteria that are more general than that of Friedman and Rubin (1967) and more parsimonious than that of Scott and Symons (1971) for the unequal-covariance case. The reparameterization of covariance matrices in terms of the eigenvalue decomposition has also been considered by Flury (1988) although he did not view it in the context of cluster analysis and he assumed the eigenvector matrices, D_k , to be the same across all groups.

A general and practical approach to non-Gaussian clustering is introduced. It is developed in detail for the important special case where points are distributed uniformly along and tightly about a line segment in p -space. "Noise" is allowed for by permitting isolated observations to be distributed over the data region according to a Poisson process. We propose an approximate Bayesian method for choosing the number of clusters. We also write down the exact Bayesian solution, which is optimal given the model, but is usually not computable; our approximation seems to perform well in numerical examples.

An alternative specification of the model (1.1), which leads to the so-called mixture maximum likelihood approach, has been considered by Wolfe (1970), Symons (1981), McLachlan (1982), and McLachlan and Basford (1988). This assumes that \mathbf{x} is a random sample from a mixture of the G densities $f_k(\mathbf{x}; \theta)$ ($k = 1, \dots, G$) in the proportions $\varepsilon = (\varepsilon_1, \dots, \varepsilon_G)^T$. Then θ and ε are estimated, and conditional probabilities $p(\gamma_i = k | \mathbf{x}, \hat{\theta}, \hat{\varepsilon})$ are calculated. Marriott (1975) and Bryant and Williamson (1978) showed that when, unlike here, estimation of θ is of primary interest, then the classification maximum likelihood method is inconsistent. However, when the covariance matrices are unequal, the mixture maximum likelihood approach appears to break down in practice (Day, 1969). McLachlan and Basford (1988, §2.1) discuss some theoretical results which suggest that it may be possible to apply the mixture maximum likelihood approach when the covariance matrices are unequal, but this does not seem to have been done yet. If it could be done, it seems likely that the methods proposed in this paper could also be extended to the mixture maximum likelihood approach using the EM algorithm (McLachlan and Basford, 1988, §1.6).

The classification and mixture maximum likelihood approaches are in conflict only when the primary aim is to estimate θ ; the conflict is resolved when, as here, the aim is to estimate γ , and θ is a nuisance parameter. This is easiest to see in a Bayesian framework, where the full solution is the

posterior distribution $p(\gamma | \mathbf{x})$. It follows from equation (2.2) of Binder (1978) that this is the same under the two models when the prior for γ in (1.1) is hierarchical and compatible with the prior for ε in the mixture model. Thus the classification maximum likelihood solution $\hat{\gamma}$ may be viewed as an approximation to the posterior mode of γ under both models.

8. Software Implementation

A comprehensive set of programs has been written by Dr Chris Fraley to implement the methods described here as well as a wide range of other clustering methods that have a model-based interpretation. The software is called “mclust” (for model-based clustering), and is available free from StatLib either as a stand-alone Fortran program, or as an S-PLUS function.

The program carries out hierarchical clustering on a data set of arbitrary dimension using any of the six criteria in Table 1, or any of the centroid, weighted average link, group average link, complete link, or single link criteria; see Gordon (1981) for definitions of these other criteria. The program allows iterative relocation if desired, and can accommodate noise explicitly with each criterion, as described in Section 3.2. It calculates the AWE for each criterion and each number of clusters.

The S-PLUS version produces output in a form that can be used by the other S-PLUS functions “plclust,” “labclust,” “cuttree,” “clorder,” and “subtree” to plot the hierarchical clustering tree, or dendrogram, with labels if desired, to create groups for further analysis within S-PLUS, and to create subtrees.

The programs may be obtained from StatLib by sending an e-mail message to “statlib@temper.stat.cmu.edu” with the single line “send mclust from general” for the Fortran version, or the single line “send mclust from S” for the S-PLUS version. While the program is not maintained, questions may be addressed by e-mail to Chris Fraley (statsci@fraley@uunet.uu.net) or to Adrian Raftery (raftery@stat.washington.edu).

ACKNOWLEDGEMENTS

Jeff Banfield’s research was supported in part by the Office of Naval Research under Contract No. N-00014-89-J-1114. Adrian Raftery’s research was supported by the Office of Naval Research under Contract No. N-00014-91-J-1074. The authors are grateful to Dr David Haynor for providing the MRI data and for helpful discussions about its medical interpretation, to Dr Chris Fraley for helpful comments and for software implementation, and to the associate editor and a referee for extremely thorough readings that led to considerable improvements in the paper.

RÉSUMÉ

L’approche par maximum de vraisemblance de la classification est suffisamment générale pour recouvrir de nombreux algorithmes classiques de regroupement, incluant ceux fondés sur le critère de la somme des carrés et sur celui de Friedman et Rubin (1967, *Journal of the American Statistical Association* **42**, 1159–1178). Néanmoins, dans son usage courant, il ne permet pas de spécifier quelles caractéristiques (orientation, taille, forme) sont communes à toutes les classes et quelles sont celles qui diffèrent entre classes.

Nous proposons des moyens pour dépasser ces limitations. Une reparamétrisation de la matrice de dispersion nous permet de spécifier que certaines caractéristiques, mais pas toutes, sont identiques pour toutes les classes. Nous précisons un cadre pratique pour des classes non-gaussiennes, et nous décrivons le moyen d’incorporer un bruit dans la forme d’un processus de Poisson. Nous proposons une méthode bayésienne approchée pour choisir le nombre de classes.

Nous étudions par simulation la performance des méthodes proposées, avec des résultats encourageants. Nous appliquons ces dernières à l’analyse de données sur le diabète, et les résultats semblent meilleurs que ceux préalablement obtenus. Nous analysons aussi une image du cerveau par résonance magnétique nucléaire, les méthodes apparaissent couronnées de succès pour faire apparaître les caractéristiques essentielles d’intérêt anatomique. Les méthodes décrites ont été codées en Fortran et en S-PLUS, elles sont disponibles librement via StatLib.

REFERENCES

- Adams, J. B., Smith, M. O., and Johnson, P. E. (1986). Spectral mixture modeling: A new analysis of rock and soil types at the Viking Lander 1 site. *Journal of Geophysical Research* **91**, 8098–8112.
- Aitkin, M., Anderson, D., and Hinde, J. (1981). Statistical modelling of data on teaching styles (with Discussion). *Journal of the Royal Statistical Society, Series A* **144**, 419–461.
- Akman, V. E. and Raftery, A. E. (1986). Bayes factors for non-homogeneous Poisson processes with vague prior information. *Journal of the Royal Statistical Society, Series B* **48**, 322–329.

- Andrews, D. F. and Herzberg, A. M. (1985). *Data: A Collection of Problems from Many Fields for Students and Research Workers*. New York: Springer-Verlag.
- Banfield, J. and Raftery, A. (1992). Ice floe identification in satellite images using mathematical morphology and clustering about principal curves. *Journal of the American Statistical Association* **87**, 7–16.
- Berger, J. O. and Sellke, T. (1987). Testing a point null hypothesis: The irreconcilability of P -values and evidence. *Journal of the American Statistical Association* **82**, 112–122.
- Binder, D. A. (1978). Bayesian cluster analysis. *Biometrika* **65**, 31–38.
- Bryant, P. and Williamson, J. A. (1978). Asymptotic behaviour of classification maximum likelihood estimates. *Biometrika* **65**, 273–278.
- Day, N. (1969). Estimating the components of a mixture of normal distributions. *Biometrika* **56**, 463–474.
- Donoho, A. W., Donoho, D. L., and Gasko, M. (1988). MACSPIN: Dynamic graphics on a desktop computer. In *Dynamic Graphics for Statistics*, W. S. Cleveland and M. E. McGill (eds), 331–352. Belmont, California: Wadsworth & Brooks/Cole.
- Everitt, B. S. (1981). Contribution to the discussion of paper by M. Aitkin, D. Anderson, and J. Hinde. *Journal of the Royal Statistical Society, Series A* **144**, 457–458.
- Flury, B. (1988). *Common Principal Components and Related Multivariate Models*. New York: Wiley.
- Friedman, H. P. and Rubin, J. (1967). On some invariant criteria for grouping data. *Journal of the American Statistical Association* **62**, 1159–1178.
- Good, I. J. (1983). *Good Thinking: The Foundation of Probability and Its Applications*. Minneapolis: University of Minnesota Press.
- Gordon, A. D. (1981). *Classification: Methods for the Exploratory Analysis of Multivariate Data*. New York: Chapman and Hall.
- Gordon, A. D. (1987). A review of hierarchical classification. *Journal of the Royal Statistical Society, Series A* **150**, 119–137.
- Jain, A. K. and Dubes, R. C. (1988). *Algorithms for Clustering Data*. Englewood Cliffs, New Jersey: Prentice-Hall.
- Maronna, R. and Jacovkis, P. M. (1974). Multivariate clustering procedures with variable metrics. *Biometrics* **30**, 499–505.
- Marriott, F. (1975). Separating mixtures of normal distributions. *Biometrics* **31**, 767–769.
- McLachlan, G. (1982). The classification and mixture maximum likelihood approaches to cluster analysis. In *Handbook of Statistics* (Vol. 2), P. R. Krishnaiah and L. N. Kanal (eds), 199–208. Amsterdam: North-Holland.
- McLachlan, G. (1987). On bootstrapping the likelihood ratio test statistic for the number of components in a normal mixture. *Applied Statistics* **36**, 318–324.
- McLachlan, G. and Basford, K. (1988). *Mixture Models: Inference and Applications to Clustering*. New York: Marcel Dekker.
- Murtagh, F. (1985). *Multidimensional Clustering Algorithms*. CompStat Lectures **4**. Heidelberg: Physica-Verlag.
- Murtagh, F. and Raftery, A. E. (1984). Fitting straight lines to point patterns. *Pattern Recognition* **17**, 479–483.
- Oldendorf, W. and Oldendorf, W. Jr. (1988). *Basics of Magnetic Resonance Imaging*. Boston: Martinus Nijhoff.
- Raftery, A. E. (1986a). A note on Bayes factors for log-linear contingency table models with vague prior information. *Journal of the Royal Statistical Society, Series B* **48**, 249–250.
- Raftery, A. E. (1986b). Choosing models for cross-classifications. *American Sociological Review* **51**, 145–146.
- Raftery, A. E. (1987). Inference and prediction for a general order statistic model with unknown population size. *Journal of the American Statistical Association* **82**, 1163–1168.
- Raftery, A. E. (1988). Analysis of a simple debugging model. *Applied Statistics* **37**, 12–22.
- Raftery, A. E. and Akman, V. E. (1986). Bayesian analysis of a change-point Poisson process. *Biometrika* **73**, 85–89.
- Reaven, G. M. and Miller, R. G. (1979). An attempt to define the nature of chemical diabetes using a multidimensional analysis. *Diabetologia* **16**, 17–24.
- Scott, A. J. and Symons, M. J. (1971). Clustering methods based on likelihood ratio criteria. *Biometrics* **27**, 387–397.
- Smith, A. and Spiegelhalter, D. (1980). Bayes factors and choice criteria for linear models. *Journal of the Royal Statistical Society, Series A* **42**, 213–220.
- Symons, M. (1981). Clustering criteria and multivariate normal mixtures. *Biometrics* **37**, 35–43.
- Vesecky, J., Samadani, R., Smith, M., Daida, J., and Bracewell, R. (1988). Observation of sea-ice dynamics using synthetic aperture radar images: Automated analysis. *IEEE Transactions in Geoscience and Remote Sensing* **GE-26**, 38–47.
- Ward, J. H. (1963). Hierarchical groupings to optimize an objective function. *Journal of the American Statistical Association* **58**, 236–244.

- Wolfe, J. H. (1970). Pattern clustering by multivariate mixture analysis. *Multivariate Behavioral Research* **5**, 329–350.
- Wolfe, J. H. (1971). A Monte Carlo study of the sampling distribution of the likelihood ratio for mixtures of multinormal distributions. Technical Bulletin STB 72-2. San Diego: U.S. Naval Personnel and Training Research Laboratory.

Received December 1989; revised November 1991 and March 1992; accepted April 1992.