

Clustering, Classification, Discriminant Analysis, and Dimension Reduction via Generalized Hyperbolic Mixtures

Katherine Morris* Paul D. McNicholas†

Abstract

A method for dimension reduction with clustering, classification, or discriminant analysis is introduced. This mixture model-based approach is based on fitting generalized hyperbolic mixtures on a reduced subspace within the paradigm of model-based clustering, classification, or discriminant analysis. A reduced subspace of the data is derived by considering the extent to which group means and group covariances vary. The members of the subspace arise through linear combinations of the original data, and are ordered by importance via the associated eigenvalues. The observations can be projected onto the subspace, resulting in a set of variables that captures most of the clustering information available. The use of generalized hyperbolic mixtures gives a robust framework capable of dealing with skewed clusters. Although dimension reduction is increasingly in demand across many application areas, the authors are most familiar with biological applications and so two of the five real data examples are within that sphere. Simulated data are also used for illustration. The approach introduced herein can be considered the most general such approach available, and so we compare results to three special and limiting cases. Comparisons with several well established techniques illustrate its promising performance.

Keywords: Dimension reduction; generalized hyperbolic distribution; mixture models; model-based clustering; model-based classification; model-based discriminant analysis.

1 Introduction

A method for estimating a projection subspace basis derived from the fit of a generalized hyperbolic mixture (HMMDR) is introduced within the paradigms of model-based clustering, classification, and discriminant analysis. This is the most general case of work in this direction over the last few years, starting with an analogous approach based on Gaussian mixtures (GMMDR; Scrucca, 2010).

Many dimension reduction methods summarize the information available through a reduced combination of the original variables. However, in terms of visualization, they do not always provide adequate information on the potential structure of the data at hand. The method proposed herein addresses this issue by revealing the underlying data clusters. At the same time, using heavy-tailed distributions, such as the generalized hyperbolic distribution, to model data can be advantageous because they assign appropriate weights to more extreme points (McNeil et al., 2005). The goal is to estimate a subspace that captures most of the clustering structure contained in the data. At the core of the method lies the sliced inverse regression (SIR) work of Li (1991, 2000), which reduces data dimensionality by considering the variation in group means to identify the subspace. Scrucca (2010) extended the SIR ideas to also include variation of group covariances. The members of the subspace arise through linear combinations of the original data, and are ordered by importance via their associated eigenvalues. The original observations in the data can be projected onto the subspace, resulting in a set of variables that captures most of the clustering information available.

The remainder of the paper is outlined as follows. Section 2 presents the background material. We then outline our dimension reduction method for selecting a reduced combination of the variables while retaining most of the clustering information contained within the data (Section 3). In Section 4, the algorithm is applied to simulated and real data

*Department of Mathematics & Statistics, University of Guelph, Ontario, Canada, N1G 2W1.

†Department of Mathematics & Statistics, McMaster University, 1280 Main St. W., Hamilton, Ontario, Canada L8S 4L8. Email: mcnicholas@math.mcmaster.ca

sets and the performance of our method is compared with its Gaussian and non-Gaussian analogues as well as with other subspace clustering techniques. Section 5 provides conclusions and suggestions for future work. Note that all computational work herein was carried out using R (R Core Team, 2013).

2 Background

2.1 Finite Mixture Models

Modern data sets used in many practical applications have grown in size and complexity, compelling the use of clustering and classification algorithms based on probability models. The model-based approach assumes that data are generated by a finite mixture of probability distributions. A p -dimensional random vector \mathbf{X} is said to arise from a parametric finite mixture distribution if its density is a convex set of probability densities, i.e.,

$$f(\mathbf{x} | \boldsymbol{\vartheta}) = \sum_{g=1}^G \pi_g f_g(\mathbf{x} | \boldsymbol{\theta}_g),$$

where G is the number of components, π_g are mixing proportions, so that $\sum_{g=1}^G \pi_g = 1$ and $\pi_g > 0$, and $\boldsymbol{\vartheta} = (\pi_1, \dots, \pi_G, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_G)$ is the parameter vector. The $f_g(\mathbf{x} | \boldsymbol{\theta}_g)$ are called component densities and $f(\mathbf{x} | \boldsymbol{\vartheta})$ is formally referred to as a G -component parametric finite mixture distribution. The use of mixture models in clustering applications can be traced back a half-century to an application of Gaussian mixture models (Wolfe, 1963). Gaussian mixture model-based approaches have been very popular due to their mathematical tractability, and until recently, they dominated literature in the field. Extensive details on finite mixture models are given by Everitt and Hand (1981), McLachlan and Basford (1988), and McLachlan and Peel (2000).

In the past several years, non-Gaussian approaches to model-based clustering, classification, and discriminant analysis have flourished. This includes work on mixtures of multivariate t -distributions (Peel and McLachlan, 2000; Greselin and Ingrassia, 2010; Andrews et al., 2011; Steane et al., 2012; Andrews and McNicholas, 2012a; McNicholas, 2013), shifted asymmetric Laplace distributions (Franczak et al., 2014), skew-normal distributions (Lin, 2010), skew t -distributions (Vrbik and McNicholas, 2012, 2014; Lee and McLachlan, 2013; Murray et al., 2014a,b), and variance-gamma distributions (McNicholas et al., 2013). Mixtures of generalized hyperbolic distributions (Browne and McNicholas, 2015) are particularly relevant to work described herein. While it is not feasible to provide an exhaustive listing here, suffice it to say that the breadth of research on non-Gaussian model-based clustering and classification is becoming as rich as that of its Gaussian precursor.

Generalized hyperbolic distributions were introduced by Barndorff-Nielsen (1977) and used to model eolian sand deposits, i.e., sand deposits arising from the action of wind. The name of the distribution was derived from the fact that its log-density has the shape of a hyperbola. Properties of generalized hyperbolic densities were discussed in Barndorff-Nielsen and Halgreen (1977) and Blæsild (1978) and, more recently, mixtures of these distributions appear in McNeil et al. (2005) and Härdle and Simar (2011). Generalized hyperbolic distributions can effectively model extreme values, making them very useful in the context of financial and risk management applications, where the normal distribution does not offer a good description of reality. The multivariate generalized hyperbolic family is extremely flexible and contains many special and limiting cases, such as the inverse Gaussian, Laplace, and skew- t distributions.

2.2 Generalized Hyperbolic Mixtures

Browne and McNicholas (2015) propose a multivariate generalized hyperbolic mixture model (HMM),

$$f(\mathbf{x} | \boldsymbol{\vartheta}) = \sum_{g=1}^G \pi_g f_h(\mathbf{x} | \lambda_g, \omega_g, \boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g, \boldsymbol{\alpha}_g), \tag{1}$$

where $\pi_g > 0$, with $\sum_{g=1}^G \pi_g = 1$, are the mixing proportions and the g th component density is

$$f_h(\mathbf{x} \mid \lambda_g, \omega_g, \boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g, \boldsymbol{\alpha}_g) = \frac{\left[\frac{\omega_g + \delta(\mathbf{x}, \boldsymbol{\mu}_g \mid \boldsymbol{\Sigma}_g)}{\omega_g + \boldsymbol{\alpha}_g^\top \boldsymbol{\Sigma}_g^{-1} \boldsymbol{\alpha}_g} \right]^{(\lambda_g - p/2)/2}}{\frac{K_{\lambda_g - p/2} \left(\sqrt{[\omega_g + \boldsymbol{\alpha}_g^\top \boldsymbol{\Sigma}_g^{-1} \boldsymbol{\alpha}_g](\omega_g + \delta(\mathbf{x}, \boldsymbol{\mu}_g \mid \boldsymbol{\Sigma}_g))} \right)}{(2\pi)^{p/2} |\boldsymbol{\Sigma}_g|^{1/2} K_{\lambda_g}(\omega_g) \exp(-(\mathbf{x} - \boldsymbol{\mu}_g)^\top \boldsymbol{\Sigma}_g^{-1} \boldsymbol{\alpha}_g)}}, \quad (2)$$

with index parameter λ_g , concentration parameter ω_g , skewness parameter $\boldsymbol{\alpha}_g$, location $\boldsymbol{\mu}_g$, and scale matrix $\boldsymbol{\Sigma}_g$. Here, $\delta(\mathbf{x}, \boldsymbol{\mu}_g \mid \boldsymbol{\Sigma}_g) = (\mathbf{x} - \boldsymbol{\mu}_g)^\top \boldsymbol{\Sigma}_g^{-1} (\mathbf{x} - \boldsymbol{\mu}_g)$ is the squared Mahalanobis distance between \mathbf{x} and $\boldsymbol{\mu}_g$ and K_{λ_g} denotes the modified Bessel function of the third kind with index λ_g .

The evaluation of modified Bessel functions in the density (2) sometimes leads to numerical overflow or underflow. To avoid these issues, we use asymptotic expansions from Abramowitz and Stegun (1972), i.e., for large x or λ ,

$$K_\lambda(\lambda x) = \sqrt{\frac{\pi}{2\lambda}} \frac{\exp\{-\lambda\rho\}}{(1+x^2)^{1/4}} \left[1 + \sum_{k=1}^{\infty} (-1)^k \frac{u_k(\tau)}{\lambda^k} \right],$$

where

$$\rho = \sqrt{1+x^2} + \ln \left(\frac{x}{1+\sqrt{1+x^2}} \right),$$

and $u_k(\tau)$ is the Debye polynomial represented by $u_0(\tau) = 1$ and

$$u_{k+1}(\tau) = \frac{1}{2}\tau^2(1-\tau^2)u'_k(\tau) + \frac{1}{8}\int_0^\tau (1-5s^2)u_k(s)ds,$$

for $\tau = 1/\sqrt{1+x^2}$ and $k = 1, 2, \dots$

The parametrization in (2) is one of several available for generalized hyperbolic distributions (cf. McNeil et al., 2005). In this case, the p -dimensional random vector \mathbf{X} is generated by combining a generalized inverse Gaussian (GIG) random variable Y with a latent multivariate Gaussian random variable $\mathbf{U} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$. Note that the density of $Y \sim \text{GIG}(\omega, \eta, \lambda)$ is

$$h(y \mid \omega, \eta, \lambda) = \frac{(y/\eta)^{\lambda-1}}{2\eta K_\lambda(\omega)} \exp \left\{ -\frac{\omega}{2} \left(\frac{y}{\eta} + \frac{\eta}{y} \right) \right\}.$$

We fix $\eta = 1$ and use the relationship $\mathbf{x} = \boldsymbol{\mu} + Y\boldsymbol{\alpha} + \sqrt{Y}\mathbf{U}$. Full details on the derivation of this parametrization and its use in parameter estimation are given by Browne and McNicholas (2015).

In the following sections, we discuss using the generalized hyperbolic distribution for model-based methods in the context of unsupervised (clustering), semi-supervised (classification) and supervised (discriminant analysis) learning. Figure 1 shows the relationship between these learning approaches.

2.3 Model-Based Clustering

Consider a clustering scenario in which none of the observations have known component membership, i.e., where the observations are unlabelled. The generalized hyperbolic model-based clustering likelihood is

$$\mathcal{L}(\boldsymbol{\vartheta} \mid \mathbf{x}) = \prod_{i=1}^n \sum_{g=1}^G \pi_g f_h(\mathbf{x}_i \mid \lambda_g, \omega_g, \boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g, \boldsymbol{\alpha}_g). \quad (3)$$

To facilitate discussion of parameter estimation, introduce Z_{ig} to denote component membership labels, so that $z_{ig} = 1$ if observation \mathbf{x}_i belongs to component g and $z_{ig} = 0$ otherwise.

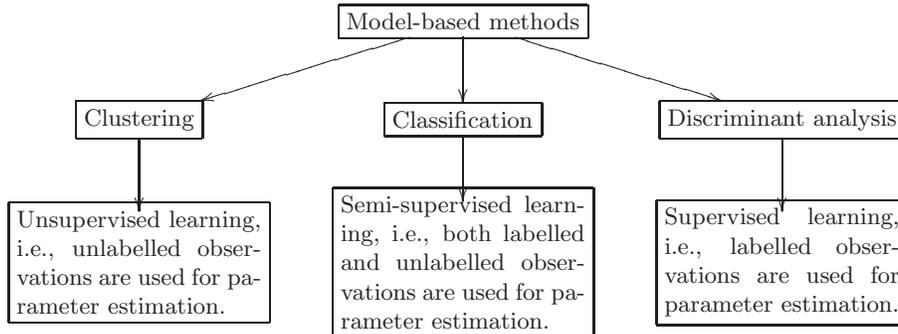


Figure 1: Diagram showing the relationship between the model-based learning paradigms discussed.

Parameter estimation for generalized hyperbolic mixtures is carried out using the expectation-maximization (EM) algorithm (Baum et al., 1970; Orchard and Woodbury, 1972; Sundberg, 1974; Dempster et al., 1977). The EM algorithm is an iterative procedure for finding maximum likelihood estimates when data are incomplete or are treated as being incomplete. The EM algorithm is based on the complete-data log-likelihood (4), where the complete-data comprise the observed \mathbf{x}_i , the missing z_{ig} , and the latent y_{ig} , for $i = 1, \dots, n$ and $g = 1, \dots, G$. Our complete-data log-likelihood can be written

$$l(\boldsymbol{\vartheta}) = \sum_{i=1}^n \sum_{g=1}^G z_{ig} \left[\log(\pi_g) + \sum_{j=1}^p \log[\phi(\mathbf{x}_i | \boldsymbol{\mu}_g + y_{ig} \boldsymbol{\alpha}_g, y_{ig} \boldsymbol{\Sigma}_g)] + \log[h(y_{ig} | \omega_g, \lambda_g)] \right], \quad (4)$$

where $\mathbf{X}_i | y_{ig} \sim \mathcal{N}(\boldsymbol{\mu}_g + y_{ig} \boldsymbol{\alpha}_g, y_{ig} \boldsymbol{\Sigma}_g)$ and $Y_{ig} \sim \text{GIG}(\omega_g, 1, \lambda_g)$. Two steps are iterated until convergence is reached. In the expectation step (E-step), the expected value of the complete-data log-likelihood is computed. Then in the maximization step (M-step), the expected value of the complete-data log-likelihood is maximized with respect to the model parameters. Extensive details of the EM algorithm for generalized hyperbolic mixtures are given by Browne and McNicholas (2015).

Following Böhning et al. (1994) and Lindsay (1995), convergence can be determined based on an asymptotic estimate of the log-likelihood at iteration $k + 1$, namely

$$l_{\infty}^{(k+1)} = l^k + \frac{l^{(k+1)} - l^{(k)}}{1 - a^{(k)}},$$

where

$$a^{(k)} = \frac{l^{(k+1)} - l^{(k)}}{l^{(k)} - l^{(k-1)}}$$

denotes the Aitken acceleration (Aitken, 1926) at iteration k . The algorithm can be considered to have converged when $l_{\infty}^{(k+1)} - l^{(k)} < \epsilon$, provided that this difference is positive (cf. McNicholas et al., 2010).

In our applications (Section 4), we assume that the number of components G is unknown. This is not unusual in real model-based clustering applications, where a criterion is often used to determine G . The Bayesian information criterion (BIC; Schwarz, 1978) is the most popular choice and is given by

$$\text{BIC} = 2l(\mathbf{x}, \hat{\boldsymbol{\vartheta}}) - r \log n,$$

where $l(\mathbf{x}, \hat{\boldsymbol{\vartheta}})$ is the maximized (observed) log-likelihood, $\hat{\boldsymbol{\vartheta}}$ denotes the maximum likelihood estimate of $\boldsymbol{\vartheta}$, r represents the number of free parameters, and n is the number of observations. After convergence, component memberships are usually estimated based on the maximum *a posteriori* (MAP) classification given by $\text{MAP}\{\hat{z}_{ig}\} = 1$ if $\arg \max_h \{\hat{z}_{ih}\} = g$ and $\text{MAP}\{\hat{z}_{ig}\} = 0$ otherwise, for $i = 1, \dots, n$.

2.4 Model-Based Classification and Discriminant Analysis

Model-based classification, or partial classification (cf. McLachlan, 1982), is a semi-supervised analogue of model-based clustering that has historically received much less attention within the literature. However, model-based classification has garnered increased attention over the past few years and some authors (e.g., Dean et al., 2006; McNicholas, 2010) have demonstrated that model-based classification can give excellent performance in real applications. Model-based discriminant analysis (Hastie and Tibshirani, 1996) is a supervised analogue of model-based clustering that has similarly received much less attention until recently (e.g., Andrews and McNicholas, 2011, 2012a). Fraley and Raftery (2002) discuss their own discriminant analysis approach, i.e., MclustDA, as well as the EDDA approach of Bensmail and Celeux (1996).

Consider the classification scenario where there are n observations, k of which have known labels. Under a model-based classification framework, all n observations are used to estimate the group memberships for the $n - k$ unlabelled observations. Following McNicholas (2010), without loss of generality, we arrange the data so that the first k observations are labelled. Accordingly, the likelihood can be written

$$\mathcal{L}(\boldsymbol{\vartheta} \mid \mathbf{x}) = \prod_{i=1}^k \prod_{g=1}^G [\pi_g f_h(\mathbf{x}_i \mid \lambda_g, \omega_g, \boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g, \boldsymbol{\alpha}_g)]^{z_{ig}} \prod_{j=k+1}^n \sum_{s=1}^G \pi_s f_h(\mathbf{x}_j \mid \lambda_s, \omega_s, \boldsymbol{\mu}_s, \boldsymbol{\Sigma}_s, \boldsymbol{\alpha}_s). \quad (5)$$

As in the case of model-based clustering, parameter estimation is carried out using the EM algorithm. From (3) and (5), we see that model-based clustering can be viewed as a special case of model-based classification that arises by considering (5) with $k = 0$.

For model-based discriminant analysis, we again have n observations, k of which have known labels. Again, we arrange the data so that the first k observations have known labels; however, instead of using all n observations to estimate the unknown labels, we only use the first k observations (i.e., the labelled observations). First, we form the likelihood

$$\mathcal{L}(\boldsymbol{\vartheta} \mid \mathbf{x}) = \prod_{i=1}^k \prod_{g=1}^G [\pi_g f_h(\mathbf{x}_i \mid \lambda_g, \omega_g, \boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g, \boldsymbol{\alpha}_g)]^{z_{ig}}. \quad (6)$$

Then, the parameter estimates are computed via the EM algorithm. The resulting *a posteriori* expected values of the Z_{ig} are used to estimate the membership labels of the remaining $n - k$ observations. In their discriminant analysis approach, Hastie and Tibshirani (1996) allow multiple Gaussian mixture components per class. Scrucca (2013) extended the dimension reduction and clustering approach of Scrucca (2010), cf. Section 2.5), to a discriminant analysis framework. However, for the discriminant analyses herein, we restrict our HMMDR approach (cf. Section 3) to one component per known class. In part, this is done because of the large number of parameters to be estimated and the relatively small number of observations in the real data sets we consider (cf. Section 4.3). However, it is also done because we believe that the flexibility inherent in generalized hyperbolic components makes it far less likely, relative to their Gaussian components, that multiple components would be needed to model a single class. This latter point will be investigated as part of future work (cf. Section 5).

2.5 Dimension Reduction and Model-Based Clustering (GMMDR)

Scrucca (2010) proposed a method of dimension reduction for model-based clustering within the Gaussian mixture framework, called GMMDR. Given a G -component Gaussian mixture model (GMM), i.e.,

$$f(\mathbf{x} \mid \boldsymbol{\vartheta}) = \sum_{g=1}^G \pi_g f_g(\mathbf{x} \mid \boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g) = \sum_{g=1}^G \pi_g \left[\frac{\exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_g)^\top \boldsymbol{\Sigma}_g^{-1} (\mathbf{x} - \boldsymbol{\mu}_g) \right\}}{(2\pi)^{\frac{p}{2}} |\boldsymbol{\Sigma}_g|^{\frac{1}{2}}} \right],$$

the procedure finds the smallest subspace that captures the clustering information contained within the data. The core of the method is to identify those directions where the cluster means $\boldsymbol{\mu}_g$ and the cluster covariances $\boldsymbol{\Sigma}_g$ vary as much as possible, provided that each direction is $\boldsymbol{\Sigma}$ -orthogonal to the others.

Finding these directions is achieved through the generalized eigen-decomposition of the kernel matrix \mathbf{M} , defined by Scrucca (2010) as $\mathbf{M}\mathbf{v}_i = l_i\boldsymbol{\Sigma}\mathbf{v}_i$, where $l_1 \geq l_2 \geq \dots \geq l_d > 0$ and $\mathbf{v}_i^\top \boldsymbol{\Sigma}\mathbf{v}_j = 1$ if $i = j$ and $\mathbf{v}_i^\top \boldsymbol{\Sigma}\mathbf{v}_j = 0$ otherwise. Note that there are $d \leq p$ directions that span the subspace. This kernel contains the variations in cluster means

$$\mathbf{M}_I = \sum_{g=1}^G \pi_g (\boldsymbol{\mu}_g - \boldsymbol{\mu})(\boldsymbol{\mu}_g - \boldsymbol{\mu})^\top$$

and variations in cluster covariances

$$\mathbf{M}_{II} = \sum_{g=1}^G \pi_g (\boldsymbol{\Sigma}_g - \bar{\boldsymbol{\Sigma}})\boldsymbol{\Sigma}^{-1}(\boldsymbol{\Sigma}_g - \bar{\boldsymbol{\Sigma}})^\top,$$

such that $\mathbf{M} = \mathbf{M}_I\boldsymbol{\Sigma}^{-1}\mathbf{M}_I + \mathbf{M}_{II}$.

Here, $\boldsymbol{\mu} = \sum_{g=1}^G \pi_g \boldsymbol{\mu}_g$ is the global mean, $\boldsymbol{\Sigma} = (1/n) \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top$ is the covariance matrix, and $\bar{\boldsymbol{\Sigma}} = \sum_{g=1}^G \pi_g \boldsymbol{\Sigma}_g$ is the pooled within-cluster covariance matrix. Parameter estimation is carried out using the `mclust` package (cf. Fraley and Raftery, 1999) for R. The `mclust` software fits a family of ten Gaussian mixture models, which is a subset of the 14 Gaussian parsimonious clustering models (GPCMs) introduced by Celeux and Govaert (1995). The GPCM family arises from the imposition of various constraints on eigen-decomposed component covariance matrices (cf. Banfield and Raftery, 1993; Celeux and Govaert, 1995; Fraley and Raftery, 2002).

2.6 Non-Gaussian Extensions to GMMDR

The work of Scrucca (2010) has already been extended to two non-Gaussian mixture settings. In the context of model-based clustering, Morris et al. (2013) proposed a t -distribution analogue of GMMDR, called t MMDR. This approach uses the t EIGEN family of models (Andrews and McNicholas, 2012a), which is a t -analogue of the GPCM family of models. The most general, i.e., unconstrained, member of the t EIGEN family is a mixture model with component density

$$f_t(\mathbf{x} \mid \boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g, \nu_g) = \frac{\Gamma(\frac{\nu_g+p}{2})|\boldsymbol{\Sigma}_g|^{-\frac{1}{2}}}{(\pi\nu_g)^{\frac{p}{2}}\Gamma(\frac{\nu_g}{2})(1 + \frac{\delta(\mathbf{x}, \boldsymbol{\mu}_g \mid \boldsymbol{\Sigma}_g)}{\nu_g})^{\frac{\nu_g+p}{2}}}, \quad (7)$$

where $\boldsymbol{\mu}_g$ is the mean, $\boldsymbol{\Sigma}_g$ is a scale matrix, ν_g is the number of degrees of freedom, and $\delta(\mathbf{x}, \boldsymbol{\mu}_g \mid \boldsymbol{\Sigma}_g)$ is defined as before. Of course, a mixture model with component density (7) is just a mixture of multivariate t -distributions, which has been applied for clustering for some time (McLachlan and Peel, 1998; Peel and McLachlan, 2000).

Morris and McNicholas (2013) developed an analogue of GMMDR for shifted asymmetric Laplace (SAL) mixtures (Franczak et al., 2014), named SALMMDR. For SAL mixtures, the component density is

$$f_s(\mathbf{x} \mid \boldsymbol{\alpha}_g, \boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g) = \frac{2 \exp\{(\mathbf{x} - \boldsymbol{\mu}_g)^\top \boldsymbol{\Sigma}_g^{-1} \boldsymbol{\alpha}_g\}}{(2\pi)^{p/2} |\boldsymbol{\Sigma}_g|^{1/2}} \left(\frac{\delta(\mathbf{x}, \boldsymbol{\mu}_g \mid \boldsymbol{\Sigma}_g)}{2 + \boldsymbol{\alpha}_g^\top \boldsymbol{\Sigma}_g^{-1} \boldsymbol{\alpha}_g} \right)^{\nu/2} K_\nu(u),$$

with mean $\boldsymbol{\mu}_g$, scale matrix $\boldsymbol{\Sigma}_g$, and skewness $\boldsymbol{\alpha}_g$. Here, $u = \sqrt{(2 + \boldsymbol{\alpha}_g^\top \boldsymbol{\Sigma}_g^{-1} \boldsymbol{\alpha}_g) \delta(\mathbf{x}, \boldsymbol{\mu}_g \mid \boldsymbol{\Sigma}_g)}$, K_ν is the modified Bessel function of the third kind with index $\nu = (2 - p)/2$, and $\delta(\mathbf{x}, \boldsymbol{\mu}_g \mid \boldsymbol{\Sigma}_g)$ is as defined previously. The use of SAL mixtures is effective for clustering data with asymmetric components, and they can perform better than Gaussian mixtures in these cases (cf. Franczak et al., 2014).

The approach introduced herein (Section 3) aims to combine the robustness offered by the t MMDR approach with the elegance and asymmetry afforded by SALMMDR.

3 Methodology

The dimension reduction approach of Scrucca (2010) is extended through development of a generalized hyperbolic analogue. We will also develop methods for model-based classification and discriminant analysis for GMMDR and all

of the non-Gaussian analogues considered herein. Recently, Scrucca (2013) also extended GMMDR to model-based discriminant analysis.

Given a generalized hyperbolic mixture (2), we wish to find a subspace $\mathcal{S}(\beta)$ where the cluster means and cluster covariances vary the most. Although μ_g is a mean and Σ_g is a covariance matrix in (2), note that they are not the mean and covariance matrix of the random variable \mathbf{X} with the density in (2), except for the special case where $\alpha_g = \mathbf{0}$. The mean of \mathbf{X} in (2) is $\tilde{\mu}_g := \mu_g + \alpha_g$, and the covariance of \mathbf{X} is $\tilde{\Sigma}_g := \Sigma_g + \alpha_g \alpha_g^\top$. Thus, we define the kernel matrix \mathbf{M}_{HMM} for generalized hyperbolic mixtures to be

$$\mathbf{M}_{\text{HMM}} = \sum_{g=1}^G \pi_g (\tilde{\mu}_g - \mu) (\tilde{\mu}_g - \mu)^\top \Sigma^{-1} \sum_{g=1}^G \pi_g (\tilde{\mu}_g - \mu) (\tilde{\mu}_g - \mu)^\top + \sum_{g=1}^G \pi_g (\tilde{\Sigma}_g - \bar{\Sigma}) \Sigma^{-1} (\tilde{\Sigma}_g - \bar{\Sigma})^\top,$$

where $\Sigma = (1/n) \sum_{i=1}^n (\mathbf{x}_i - \mu)(\mathbf{x}_i - \mu)^\top$ denotes the overall covariance matrix and $\bar{\Sigma} = \sum_{g=1}^G \pi_g \tilde{\Sigma}_g$ is the pooled within-cluster covariance matrix.

Proposition 3.1 *The directions where the cluster means $\tilde{\mu}_g$ and the cluster covariances $\tilde{\Sigma}_g$ vary the most are obtained from the eigen-decomposition*

$$\mathbf{M}_{\text{HMM}} \mathbf{v}_i = l_i \Sigma \mathbf{v}_i, \quad (8)$$

where $l_1 \geq l_2 \geq \dots \geq l_d > 0$ and $\mathbf{v}_i^\top \Sigma \mathbf{v}_j = 1$ if $i = j$ and $\mathbf{v}_i^\top \Sigma \mathbf{v}_j = 0$ otherwise.

The eigenvectors $[\mathbf{v}_1, \dots, \mathbf{v}_d] \equiv \beta$, with $d \leq p$, form the basis of the dimension reduction subspace $\mathcal{S}(\beta)$. These eigenvectors are defined as the HMMDR directions.

Proposition 3.2 *Let $\mathcal{S}(\beta)$ be the subspace spanned by the HMMDR directions obtained from the eigen-decomposition of \mathbf{M}_{HMM} (8).*

i The projections of the parameters onto $\mathcal{S}(\beta)$ are given by $\beta^\top \tilde{\mu}_g$ and $\beta^\top \tilde{\Sigma}_g \beta$, respectively.

ii The projections of the $n \times p$ data matrix \mathbf{x} onto the subspace $\mathcal{S}(\beta)$ are computed from $\mathbf{x}\beta$. These projections are defined as the HMMDR variables.

For an $n \times p$ data matrix \mathbf{x} , the kernel \mathbf{M}_{HMM} (8) is obtained using the estimates from the fit of an HMM on \mathbf{x} , via an EM algorithm. Then, the HMMDR directions are calculated from the generalized eigen-decomposition of \mathbf{M}_{HMM} (8) with respect to the overall covariance matrix Σ . The HMMDR directions are ordered based on eigenvalues, which means that directions associated with eigenvalues close to zero can be disregarded in practical applications because clusters will superimpose greatly along these directions.

Similar to GMMDR, the estimation of the HMMDR variables can be interpreted as feature selection, where the members are reduced through a set of linear combinations of the original variables. It is possible that this set of features contains estimated HMMDR variables that do not offer any clustering information but require parameter estimation. Scrucca (2010) uses the selection method of Raftery and Dean (2006) to prune the subset of GMMDR features. We follow this approach to select the most appropriate HMMDR variables. Two subsets of features, s and $s' = \{s \setminus i\} \subset s$, for example, can be compared using the BIC difference

$$\text{BIC}_{\text{diff}}(Z_{i \in s}) = \text{BIC}_{\text{clust}}(Z_s) - \text{BIC}_{\text{not clust}}(Z_s) = \underbrace{\text{BIC}_{\text{clust}}(Z_s)}_1 - [\underbrace{\text{BIC}_{\text{clust}}(Z_{s'})}_2 + \underbrace{\text{BIC}_{\text{reg}}(Z_i | Z_{s'})}_3], \quad (9)$$

where term 1 in (9) denotes the BIC value for the best clustering model fitted using features in s , term 2 denotes the BIC value for the best clustering model fitted using features in s' , and term 3 denotes the BIC value for the regression of the i th feature on the remaining features in s' .

Because the space of all possible subsets contains $2^d - 1$ elements, where $d \leq p$, a full feature search is not usually feasible. To this end, we employ the forward greedy search algorithm of Scrucca (2010) to find a local optimum in the model space. The procedure is based on the forward-backward search algorithm of Raftery and Dean (2006); however, a backward step is not necessary here because the HMMDR variables are Σ -orthogonal. The method can be summarized

into three main stages. The initial step selects the first feature that maximizes the BIC difference in (9), between the best clustering model and the model that assumes no clustering, i.e., a single component. The following step selects the next feature amongst those not previously included to be the one that maximizes the BIC difference in (9). This process is iterated until all of the BIC differences for the inclusion of a variable become negative.

At each stage, the search over the model space is performed with respect to the model parameterization and the number of clusters. This algorithm can be applied to the three frameworks under consideration: model-based clustering, classification, and discriminant analysis, by modifying the likelihood functions (3), (5), and (6), respectively, via the EM procedure. We can now summarize our new method, which we call HMMDR:

1. Fit an HMM (1) to the data using the EM algorithm.
2. Estimate the HMMDR directions: identify directions where the cluster means and cluster variances vary the most, provided each direction is Σ -orthogonal to the others. This is done through the eigen-decomposition of the kernel matrix M_{HMM} in (8).
3. Select the HMMDR variables: compute the set of features by projecting the data onto the estimated subspace and use the greedy search algorithm to discard the ones that provide no clustering information.
4. Fit an HMM (1) on the selected HMMDR variables and return to step 2.
5. Repeat steps 2–4 until none of the features can be discarded.

Note that in the analyses herein (Section 4) the `hclust()` from the `mclust` package (Fraley et al., 2012) for R is used for initialization of the EM algorithm in step 1.

4 Applications

4.1 Performance assessment

Although the examples herein are treated as genuine clustering and classification analyses, we know the true class labels in all cases. Therefore, we can compare our predicted classifications to the true class labels in each case. To do this, we use the adjusted Rand index (ARI; Hubert and Arabie, 1985), which is the Rand index (Rand, 1971) corrected for chance agreement. The Rand index is based on pairwise agreement and takes a value between 0 and 1, where 1 indicates perfect agreement between two partitions. The correction that leads to the ARI accounts for the fact that random classification is expected to result in some correct agreements; accordingly, the ARI has an expected value of 0 under random classification and, as with the Rand index, perfect classification corresponds to a value of 1. Negative ARI values are possible and indicate classification results worse than would be expected under random classification.

4.2 Simulated data

First, we employ a data simulation scheme based on two scenarios to test the HMMDR algorithm and compare it to its Gaussian analogue. In Scenario I, we generate three variables from a mixture of multivariate Gaussian distributions with mixing proportions $\pi_1 = \pi_2 = \pi_3 = 1/3$, means $\mu_1 = (0, -2, 0)'$, $\mu_2 = (2, 4, 0)'$, $\mu_3 = (-2, -4, 2)'$, common covariance matrix $\Sigma = 0.5\mathbf{I}_3$, where \mathbf{I}_3 is the 3×3 identity matrix, and three different sample sizes $n \in \{100, 500, 1000\}$ (e.g., Figure 2). In Scenario II, we modify Scenario I by adding five noise variables generated from standard normal distributions (e.g., Figure 3).

We run each scenario 300 times for each possible combination of data dimension and model framework. For model-based classification and discriminant analysis, we assumed that each observation had a 50% probability of being known. This resulted in the number of known observations k being close to $n/2$ but varying slightly from run to run (more details on the selection of known observations appear in the next section). The ARI values for both model-based classification and discriminant analysis was computed based only on the unlabelled observations. The results (Tables 1 and 2) show that HMMDR generally exhibits high ARI values across both scenarios and all sample sizes for model-based clustering, classification, and discriminant analysis. We notice a slight drop in classification performance for Scenario II with 100

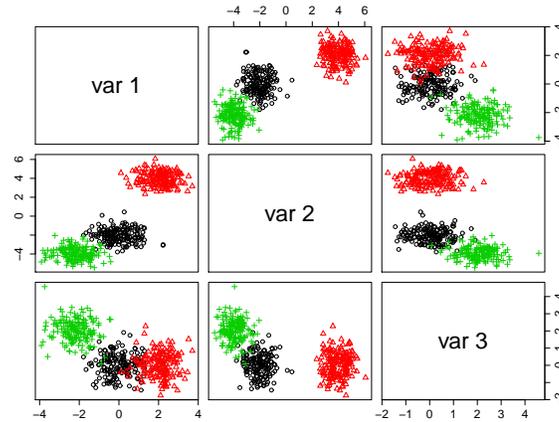


Figure 2: Pairs plot illustrating a generated data set from Scenario I.

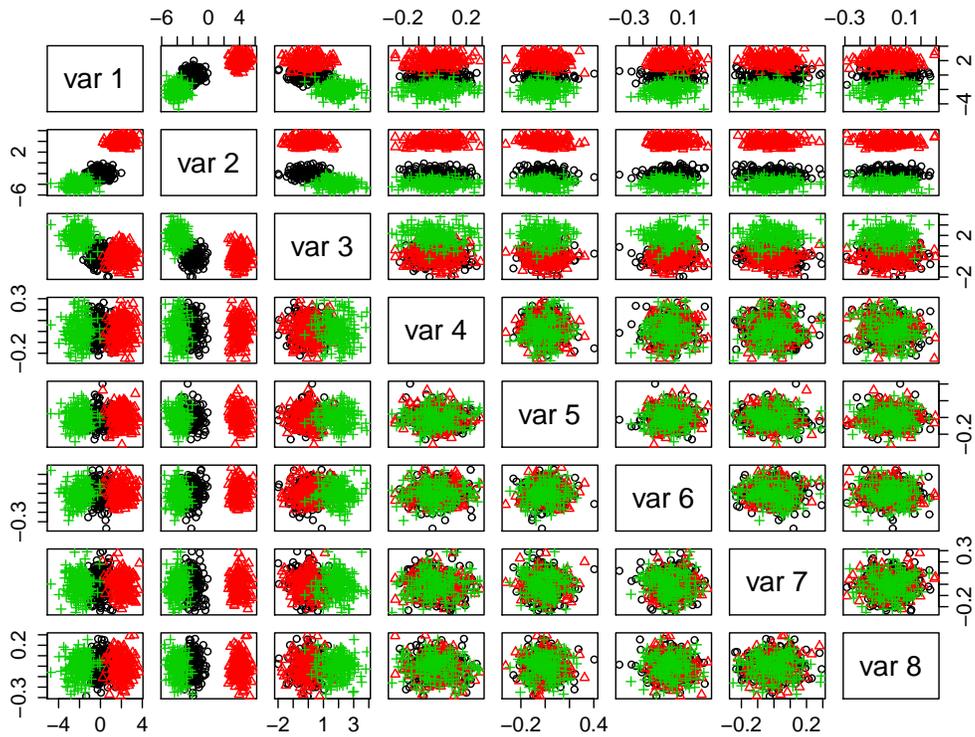


Figure 3: Pairs plot illustrating a generated data set from Scenario II.

observations; however, performance on noisy data improves for larger values of n . One or two features are selected throughout Scenario I, and, as expected, the addition of noise sometimes leads to a slight increase in the number of features in Scenario II.

The results in Tables 1 and 2 demonstrate that the HMMDR approach gives excellent performance — for model-based clustering, classification, and discriminant analysis — for the data from Scenarios I and II. However, it is also helpful to have a sense of how long the algorithms take to run. Consider the cases with $n = 1000$. For clustering, one run of the HMMDR approach takes an average of 18.5 s in Scenario I and 37.7 s in Scenario II. For classification, the equivalent times are 15.4 s and 26.1 s, and for discriminant analysis, the times are 15.0 s and 31.3 s. While it is true

Table 1: Summary of results for the HMMDR and GMMDR approaches on the simulated data from Scenario I, based on 300 runs.

		HMMDR			GMMDR		
		$n = 100$	$n = 500$	$n = 1000$	$n = 100$	$n = 500$	$n = 1000$
Clustering	Avg. ARI	0.9701	0.9634	0.9634	0.9637	0.9663	0.9668
	Std. dev. ARI	0.0349	0.0173	0.0125	0.0347	0.0166	0.0111
	Features	1-2	1-2	1	1	1	1
	Avg. no. features	1.137	1.003	1	1	1	1
Classification	Avg. ARI	0.9877	0.9829	0.9762	0.9639	0.9839	0.9872
	Std. dev. ARI	0.0223	0.0123	0.0103	0.0433	0.0010	0.0008
	Features	1-2	1-2	1-2	1	1	1
	Avg. no. features	1.07	1.033	1.02	1	1	1
Discriminant analysis	Avg. ARI	0.9274	0.9815	0.9744	0.8742	0.9656	0.9602
	Std. dev. ARI	0.1031	0.0181	0.0193	0.1108	0.046	0.038
	Features	1-2	1-2	1-2	1-2	1	1
	Avg. no. features	1.303	1.027	1.01	1.11	1	1

Table 2: Summary of results for the HMMDR and GMMDR approaches on the simulated data from Scenario II, based on 300 runs.

		HMMDR			GMMDR		
		$n = 100$	$n = 500$	$n = 1000$	$n = 100$	$n = 500$	$n = 1000$
Clustering	Avg. ARI	0.9239	0.9810	0.9827	0.9632	0.9655	0.9656
	Std. dev. ARI	0.0692	0.0120	0.0072	0.0396	0.0149	0.0105
	Features	1-4	1-3	1-3	1,4	1	1
	Avg. no. features	2.609	2.287	2.277	1.021	1	1
Classification	Avg. ARI	0.9507	0.9296	0.9616	0.9644	0.9829	0.9830
	Std. dev. ARI	0.0533	0.0129	0.0078	0.0472	0.0192	0.0072
	Features	1-4	1-3	1-3	1,4	1	1
	Avg. no. features	1.07	1.65	1.74	1.01	3	3
Discriminant analysis	Avg. ARI	0.7299	0.9662	0.9877	0.8336	0.9634	0.9609
	Std. dev. ARI	0.1280	0.0247	0.0072	0.0978	0.0354	0.0316
	Features	1-3	1-3	1-3	1-3	1-3	1-2
	Avg. no. features	1.557	1.97	1.797	2.425	1.66	1.137

that our approach would be considered slow when compared some competitors, these times show that the algorithm does not take a prohibitively long time to run.

In Scenario III, we compare the performance of HMMDR against the existing GMMDR approach for higher dimensional data, including the case where $n_g < p$. We generated three-component data sets with $n_g = 40$ observations per component from generalized hyperbolic distributions with random covariance matrices (produced using the R package `clusterGeneration`; Qiu and Joe, 2006) and $\alpha = -\mathbf{1}$. The means were drawn from a multivariate standard normal distribution and multiplied by a small integer. A typical data set from Scenario III is illustrated in Figure 4, and clearly this is a difficult clustering problem. Runs were performed with $p \in \{20, 30, 40, 50\}$, most completed successfully (i.e., converged) but some did not. Looking at the results (Table 3), we see that HMMDR tends to produce more successful runs than GMMDR when $p \leq n_g$. We note numerical difficulties for $p > n_g$, where the HMMDR scale matrices will generally be numerically singular; this problem is mitigated in GMMDR because the parsimonious covariance structures from `mclust` are used. When successful, HMMDR performs well in higher dimensions ($p = 40, 50$) but requires more features than GMMDR.

We note that Scenario III is not a fair comparison because GMMDR has the massive parsimony advantage of drawing on the `mclust` covariance structures. Repeating the simulations under this scenario with GMMDR restricted to the “VVV” model, i.e., without any reduction in the number of free covariance parameters, would tell a different story. Of course, implementing the GPCM covariance structures on the scale matrices within the HMMDR approach would also

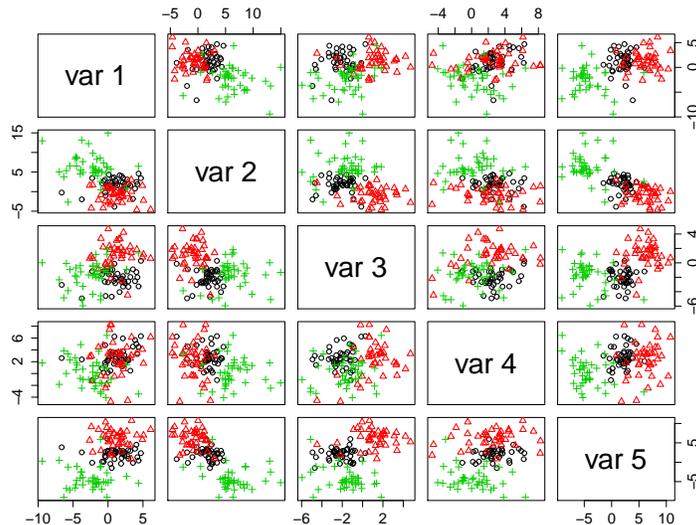


Figure 4: Pairs plot illustrating a generated data set from Scenario III.

Table 3: Summary of clustering results for the HMMDR and GMMDR approaches on the simulated data from Scenario III.

	p	Runs	Succ.	Avg. ARI	Med. ARI	Avg. Feat.	Avg. Comp.
HMMDR	20	30	30	0.43	0.41	9.8	2.93
	30	30	29	0.72	0.88	19	3
	40	50	44	0.94	0.96	25.1	3
	50	60	31	0.97	1	29.5	3
GMMDR	20	30	26	0.88	0.91	4.7	3.2
	30	30	26	0.96	1	3.2	3.1
	40	50	43	0.98	1	3.8	3
	50	60	50	0.99	1	3.7	3

lead to different results. Furthermore, we anticipate that the performance of HMMDR would also improve with the addition of other methods to avoid numerical singularities in the estimation of the scale matrices.

4.3 Real data

To gauge the performance of our algorithm on real data, we compare results with the eight methods outlined below. Except for k -means, we choose these particular comparator methods because they provide model-based analyses while implicitly reducing the dimensionality.

1. Robust principal component analysis (ROBPCA; Hubert et al., 2005) paired with t -mixtures via the t EIGEN family: principal components analysis resistant to outliers, with robust loadings computed by using projection-pursuit techniques and the minimum covariance determinant method. We use the R package `rrcov` (Todorov and Filzmoser, 2009) for the ROBPCA computations as well as the `teigen` package (Andrews and McNicholas, 2012b).
2. The family of parsimonious Gaussian mixture models (McNicholas and Murphy, 2008), which contains the mixture of factor analyzers model and variants thereof. The R package `pgmm` (McNicholas et al., 2011) is used.
3. Mixtures of common factor analyzers (Baek et al., 2010) using the the R package `mcfa` (Baek et al., 2009).

4. FisherEM (Bouveyron and Brunet, 2012): a subspace clustering method based on Gaussian mixtures, where an EM-like algorithm estimates both the discriminative subspace and the parameters of the model. The R package `FisherEM` (Bouveyron and Brunet, 2012) is employed.
5. The R package `clustvarsel` (Scrucca et al., 2013).
6. k -means clustering using the R function `kmeans`.
7. GMMDR: the approach of Scrucca (2010) based on Gaussian mixtures. While the dimension reduction procedure of GMMDR is available in the R package `mclust`, the subset selection procedure is not currently available.
8. t MMDR (Morris et al., 2013): the t -analogue of GMMDR. Fitting of the t -mixtures was carried out with the R package `teigen`.
9. SALMMDR (Morris and McNicholas, 2013): the SAL analogue of GMMDR.

For the analyses in this section, we fit HMMDR and comparator methods to the scaled version of each data set. Where appropriate, we initialize the algorithms with the Gaussian hierarchical agglomerative procedure from `mclust`. In the case of HMMDR and its analogues, we allow the number of components to vary between $G = 1$ and $G = 6$. Note that we use the term ‘analogue’ somewhat loosely here, because we do not consider decomposed covariance structures for either generalized hyperbolic mixtures or shifted asymmetric Laplace mixtures (Table 4).

Table 4: Details about the four MMDR methods.

Method	Covariance	Eigen-decomposed Σ_g	Model family
HMMDR	$\Sigma_g + \alpha_g \alpha_g^\top$	No	–
SALMMDR	$\Sigma_g + \alpha_g \alpha_g^\top$	No	–
t MMDR	$\frac{\nu_g}{\nu_g - 2} \Sigma_g, \nu_g > 2$	Yes	t EIGEN
GMMDR	Σ_g	Yes	MCLUST

In the context of model-based classification and discriminant analysis, we use the approach of McNicholas (2010) to simulate a situation in which some of the group memberships are unknown. For each observation \mathbf{x}_i , a random number is generated from a uniform distribution on $[0, 1]$. If the random number is less than 0.5, then \mathbf{x}_i is taken as known; otherwise, \mathbf{x}_i is taken as unknown. To make sure that all the classes are represented, we repeated this procedure for each group until at least one known observation was produced before moving onto the next group. Of course, it follows from this procedure that the number of unknown observations varied from run to run.

We utilize the functionality of `teigen` to fit both Gaussian mixtures and t -mixtures for model-based classification (cf. Andrews and McNicholas, 2012a). Similarly, `teigen` and `mclust` were employed for model-based discriminant analysis. For each data set, the procedures were run 25 times, using hierarchical agglomerative starting values of the unknown \hat{z}_{ig} . Note that we choose each real data set on the basis that it has previously been used to illustrate the performance of some of the comparator methods. We consider that this approach facilitates a very fair comparison. In all cases, we illustrate model-based clustering, classification, and discriminant analysis. For model-based classification and discriminant analysis, the ARI is computed based only on unlabelled observations. We perform random subset cross-validation, training on 25 different subsets consisting of roughly half the number of observations. This is more challenging than other well known procedures such as 10-fold cross-validation.

4.3.1 Swiss Bank Notes

Flury and Riedwyl (1988) present six measurements (length, diagonal, left, right, top, and bottom) taken from genuine and counterfeit Swiss bank notes. These data are available through the R package `gclus` (Hurley, 2010). In terms of model-based clustering, HMMDR and its comparators were fitted to these data and the resulting MAP classifications show very high ARI values for most methods (Table 5), with HMMDR and k -means being the only methods to cluster the bank notes data perfectly. For model-based classification, HMMDR, SALMMDR, and t MMDR produce perfect classifications of the unknown observations (Table 7). However, only HMMDR provides perfect model-based

discriminant analysis results on the bank notes. We note that HMMDR selected the minimum number of features in all three scenarios, i.e., one feature.

Table 5: Summary of model-based clustering results for the Swiss bank notes and female voles data sets.

	Bank Notes			Female Voles		
	ARI	Features	Components	ARI	Features	Components
HMMDR	1	1	2	1	1	2
SALMMDR	0.98	3	2	0.95	2	2
tMMDR	0.98	2	2	0.91	1	2
GMMDR	0.98	2	2	0.91	1	2
ROBPCA	0.98	4	2	0.91	3	2
FisherEM	0.98	1	2	0.66	1	2
clustvarsel	0.85	4	3	0.91	3	2
mcfa	0.98	2	2	0.91	2	2
pgmm	0.82	2	4	0.91	1	2
kmeans	1	–	2	0.74	–	2

It is notable that some of the clustering approaches based on a Gaussian mixture, e.g., `clustvarsel` and `pgmm`, return $G > 2$ components (Table 5). In this context, it is interesting that GMMDR returns a $G = 2$ component solution. The fact that HMMDR returns $G = 2$ components is less surprising because Tortora et al. (2015) fit four different non-Gaussian mixture approaches, including a mixture of generalized hyperbolic distributions, to the `banknote` data and choose a model with $G = 2$ components in all four cases. Using another non-Gaussian mixture approach, Franczak et al. (2015) also find a $G = 2$ component solution.

4.3.2 Female Voles

Flury (1997) discuss seven measurements (Table 6) of female voles from two species (*Microtus californicus* and *Microtus ochrogaster*) originally studied by Airoidi and Hoffmann (1984). The data are available within the R package `Flury` (Flury, 2010). Tables 5 and 7 indicate that, out of all of the procedures fitted to the voles data, HMMDR is the only one giving perfect classification results in all three paradigms.

Table 6: Measurements taken for the female vole data.

Age in days	Incisive foramen length
Condylar incisive length	Skull height
Alveolar length of upper molar tooth row	Interorbital width
Zygomatic width	

Table 7: Summary of model-based classification and discriminant analysis results for the bank note data, based on 25 runs.

	Bank Notes			Female Voles		
	ARI	Features	Components	ARI	Features	Components
HMMDR class.	1	2	2	1	1	2
SALMDR class.	1	2	2	1	3	2
<i>t</i> MDR class.	1	2	2	0.95	1	2
GMDR class.	0.96	2	2	0.96	1	2
HMMDR DA	1	2	2	1	2	2
SALMDR DA	0.96	2	2	0.91	3	2
<i>t</i> MDR DA	0.95	1–2	2	0.91	1	2
GMDR DA	0.98	1–4	2	0.88	1–4	2

4.3.3 Italian Wines

Forina et al. (1986) recorded several chemical and physical properties for three types of Italian wines: Barolo, Grignolino, and Barbera. As shown in Table 8, thirteen properties for 178 wines are available from the R package `gclus` (Hurley, 2004).

Table 8: Chemical and physical properties available for the wine data in `gclus`.

Alcohol	Proline	OD280/OD315 of diluted wines
Malic acid	Ash	Alkalinity of ash
Hue	Total phenols	Magnesium
Color intensity	Nonflavonoid phenols	Proanthocyanins
Flavonoids		

Within the model-based clustering framework, HMMDR is the best performer (ARI = 0.97, Table 10), with only two misclassified observations (Table 9). The model-based classification scenario (Table 11) reveals that HMMDR, SALMDR, and *t*MDR produce perfect classification results. With an ARI of 0.92, HMMDR gives the best performance within the model-based discriminant analysis paradigm.

Table 9: Model-based clustering, classification, and discriminant analysis results for our HMMDR approach fitted to the wine data. Model-based classification and discriminant analysis results are based on 25 runs.

	Clustering			Classification			Disc. Anal.		
	1	2	3	1	2	3	1	2	3
Barolo	59	0	0	875	0	0	625	50	0
Grignolino	0	68	2	0	925	0	0	825	0
Barbera	0	0	48	0	0	500	0	0	550
ARI; Features	0.97; 7			1; 5			0.92; 8		

Figure 5 illustrates three of the estimated HMMDR directions obtained from our model-based clustering of the wine data (Table 9). The edge histograms depict the distribution of the observations from the estimated directions, coloured by estimated cluster allocation. The plot on the left-hand side reveals quite clearly the underlying cluster structure in the data. Although there is some overlap in these two directions, only two wines were misclassified by the HMMDR method and so some of the other dimensions must give additional clarity.

Table 10: Summary of model-based clustering results for the wine data.

Method	ARI	Features	Components
HMMDR	0.97	7	3
SALMMDR	0.92	10	3
<i>t</i> MMDR	0.93	4	3
GMMDR	0.85	5	3
ROBPCA	0.83	4	3
FisherEM	0.91	2	3
clustvarsel	0.78	5	3
mcfa	0.90	3	3
pgmm	0.79	2	4
kmeans	0.90	-	3

Table 11: Summary of model-based classification and discriminant analysis results for the wine data, based on 25 runs.

Method	ARI	Features	Components
HMMDR class.	1	5	3
SALMMDR class.	1	10	3
<i>t</i> MMDR class.	1	5	3
GMMDR class.	0.93	3-8	3
HMMDR DA	0.92	8	3
SALMMDR DA	0.88	8	3
<i>t</i> MMDR DA	0.85	1-8	3
GMMDR DA	0.85	2-8	3

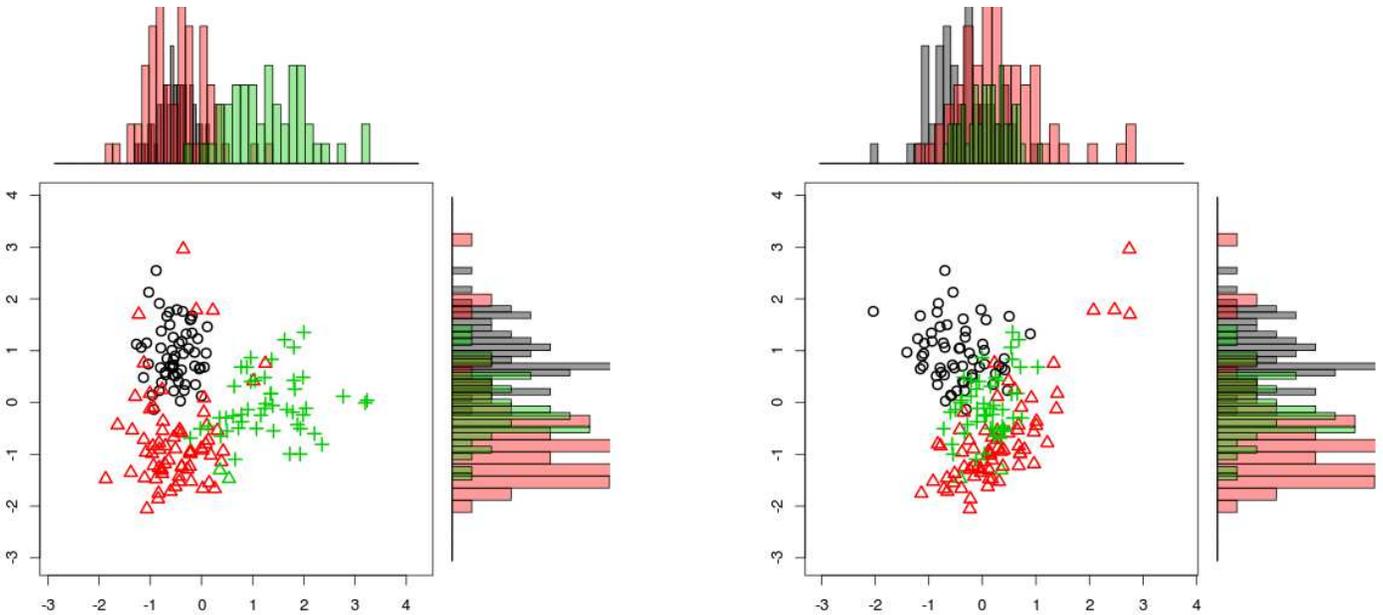


Figure 5: Plots of some of the estimated HMMDR directions for the wine data obtained from model-based clustering (directions 2 vs. 4 on the left-hand side, and directions 4 vs. 5 on the right-hand side). Symbols indicate true cluster membership and colours indicate the estimated HMMDR cluster allocation. The edge histograms depict the estimated distributions of the observations in each cluster.

4.3.4 Wisconsin Breast Cancer Study

Mangasarian et al. (1995) presented a study of breast cancer from Wisconsin, undertaken to establish whether fine needle aspiration of breast tissue samples could classify tumour status. Several attributes are recorded (Table 12) for 681 cases of potentially cancerous tumours, of which 238 were actually malignant. These data are available in the R package `faraway` (Faraway, 2011).

Within the model-based clustering framework, HMMDR selected three features and gave the best classification performance (ARI = 0.89, Table 13). The *t*MMDR (ARI = 0.86) and *k*-means (ARI = 0.84) approaches are close behind; however, the rest of the comparators did not produce particularly good results on these data (Table 14). In the model-based classification and discriminant analysis scenarios, HMMDR gave classification performance similar to

Table 12: Tissue sample properties of the Wisconsin breast cancer data.

Marginal adhesion	Epithelial cell size	Clump thickness
Bare nuclei	Mitoses	Cell shape uniformity
Bland chromatin	Normal nucleoli	Cell size uniformity

clustering and was again the best performer over the 25 runs (Table 15).

Table 13: Model-based clustering, classification, and discriminant analysis results for our HMMDR approach on the breast cancer data. Model-based classification and discriminant analysis results are based on 25 runs.

	Clustering		Classification		Disc. Anal.	
	1	2	1	2	1	2
Malignant	236	2	2925	125	2575	200
Benign	17	426	100	4925	50	4675
ARI; Features	0.89; 3		0.89; 2		0.87; 6	

Table 14: Summary of model-based clustering results for the breast cancer data.

Method	ARI	Features	Comp.
HMMDR	0.89	3	2
SALMMDR	0.80	5	2
<i>t</i> MMDR	0.86	3	2
GMMDR	0.58	5	2
ROBPCA	0.55	5	3
FisherEM	0.79	1	2
clustvarsel	0.78	2	2
mcfa	0.65	3	2
pgmm	0.42	2	4
kmeans	0.84	–	2

Table 15: Summary of model-based classification and discriminant analysis results for the breast tissue data, based on 25 runs.

Method	ARI	Features	Comp.
HMMDR class.	0.89	2	2
SALMMDR class.	0.85	4	2
<i>t</i> MMDR class.	0.86	1–2	2
GMMDR class.	0.86	1	2
HMMDR DA	0.87	6	2
SALMMDR DA	0.84	7	2
<i>t</i> MMDR DA	0.85	2	2
GMMDR DA	0.84	1–7	2

Figure 6 illustrates three of the estimated HMMDR directions obtained from the model-based clustering output shown in Table 13. The plots depict quite clearly the inherent cluster structure in the data, and we notice that the malignant breast tissues are quite tightly packed in their cluster. The red histogram on the top edge of the left-hand side of Figure 6 is a nice illustration of a skew distribution along a HMMDR direction.

4.3.5 Colon Cancer

Alon et al. (1999) analyzed gene expression data from microarray experiments of colon tissue, probed by oligonucleotide arrays. A reduced data set containing 62 tissue samples — 40 tumour tissues and 22 normal tissues — and 2,000 genes is available in the R package `plsgenomics` (Boulesteix et al., 2011). The challenge here is dealing with the large number of gene expression levels compared with the small number of tissue samples. As is the case with microarray data in general, there are many non-informative genes that can obstruct the clustering of the samples. Thus, gene filtering was carried out prior to further analysis. While McLachlan et al. (2002) and McNicholas and Murphy (2010) used the EMMIX-GENE procedure (McLachlan et al., 2002) to reduce the dimensionality of the colon data, we considered a different method.

Our gene filtering approach is to find differentially expressed genes based on modified *t*-tests, using the R Bioconductor package `siggenes` (Schwender, 2012). This package contains the function `sam`, which implements the significance analysis of microarrays (SAM) method proposed by Tusher et al. (2001). SAM computes a statistic d_i for each gene i , measuring the strength of the relationship between gene expression and the response variable (which is the class variable

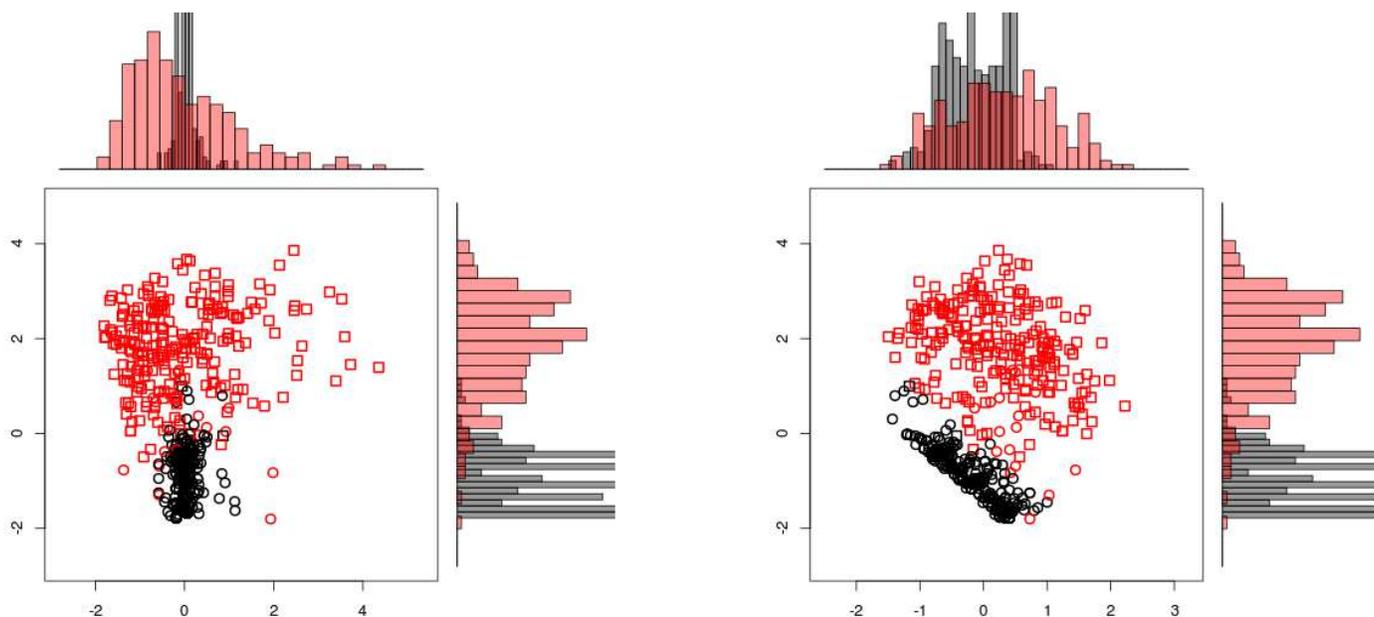


Figure 6: Plots of some of the estimated HMMDR directions for the breast cancer data obtained from model-based clustering (direction 1 vs. 3 on the left-hand side, and direction 2 vs. 3 on the right-hand side). Symbols indicate true cluster membership and colours indicate the estimated HMMDR cluster allocation. The edge histograms depict the estimated distributions of the observations in each cluster.

in our case). It uses repeated permutations of the data to determine if the expression of any gene is significantly related to the response. The cutoff for significance is determined by a tuning parameter Δ , chosen by the user based on the false positive rate. We employed 100 permutations and chose $\Delta = 2.4$, which yielded 23 genes for analysis (Table 19, A).

Even with the dimensionality reduced to 23 genes, the analysis of the colon data was quite challenging. When HMMDR was fitted to these data within the model-based clustering paradigm, five observations were misclassified (Table 16), corresponding to an ARI of 0.70. This was the best result, with t MMDR (ARI = 0.64) being the next best performer (Table 17). For model-based classification and discriminant analysis, HMMDR gave better performance with ARI values of 0.86 and 0.83, respectively (Table 18). We note that SALMMDR, t MMDR, and GMMDR also gave improved performance within the model-based classification and discriminant analysis paradigms; however, HMMDR was the best approach across all paradigms.

Table 16: Model-based clustering, classification, and discriminant analysis results for the colon data. Model-based classification and discriminant analysis results are based on 25 runs.

	Clustering		Classification		Disc. Anal.	
	1	2	1	2	1	2
Normal	22	0	225	0	220	15
Tumour	5	35	25	425	15	420
ARI; Features	0.70; 3		0.86; 5		0.83; 6	

Figure 7 illustrates three of the estimated HMMDR directions obtained from the model-based clustering output from Table 16. The plots depict the inherent cluster structure in the colon tissues and, as we would expect, the misclassified tissues are generally close to the cluster boundaries.

It is interesting to note that McNicholas and Murphy (2010) obtained similar clustering results to HMMDR for the

Table 17: Summary of model-based clustering results for the best models fitted to the colon data.

Method	ARI	Feat.	Comp.
HMMDR	0.70	3	2
SALMMDR	0.59	10	2
<i>t</i> MMDR	0.64	1	2
GMMDR	0.59	1	2
ROBPCA	0.36	3	3
FisherEM	0.59	1	2
clustvarsel	0.35	3	4
mcfa	0.64	5	2
pgmm	0.40	3	2
kmeans	0.59	-	2

Table 18: Summary of model-based classification and discriminant analysis results for the colon data, based on 25 runs.

Method	ARI	Feat.	Comp.
HMMDR class.	0.86	5	2
SALMMDR class.	0.71	6	2
<i>t</i> MMDR class.	0.75	1	2
GMMDR class.	0.73	1	2
HMMDR DA	0.83	6	2
SALMMDR DA	0.70	7	2
<i>t</i> MMDR DA	0.74	1	2
GMMDR DA	0.70	1	2

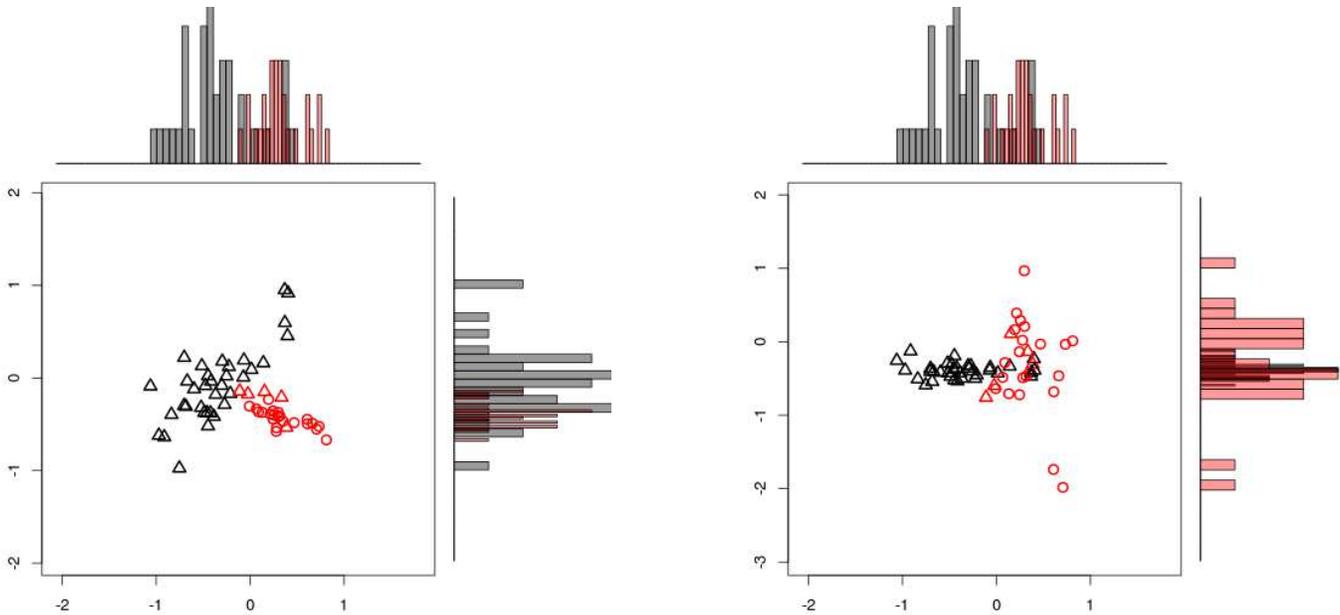


Figure 7: Plots of some of the estimated HMMDR directions for the colon data obtained from model-based clustering (direction 2 vs. 3 on the left-hand side, and direction 1 vs. 2 on the right-hand side). Symbols indicate true cluster membership and colours indicate the estimated HMMDR cluster allocation. The edge histograms depict the estimated distributions of the observations in each cluster.

colon data. They considered a subset of 461 genes and the best model fitted to these data had six latent factors, with five misclassified tissues and an ARI of 0.70. Although equal to two significant figures, we point out for completeness that our HMMDR approach has a very slightly higher ARI value (0.699 vs. 0.697). However, we must also point out that we used knowledge of tissue type in our gene selection while McNicholas and Murphy (2010) did not. McLachlan et al. (2002) also analyzed these data; they selected 446 genes and identified five types of clusterings on this subset. However, these clusterings did not correspond to the tissue type.

5 Conclusions

This paper introduced an effective dimension reduction technique for model-based clustering, classification, and discriminant analysis using multivariate mixtures of generalized hyperbolic distributions. Our method, known as HMMDR, focused on identifying the smallest subspace of the data that captured the inherent cluster structure. The HMMDR approach was illustrated using simulated and real data, where it performed favourably compared to its existing special cases, i.e., GMMDR, *t*MMDR, and SALMMDR. In clustering applications, HMMDR consistently outperformed several other model-based dimension reduction methods (ROBPCA, *pgmm*, *FisherEM*, *clustvarsel*, and *mcfa*).

One limitation sometimes encountered using our current approach is singularities while fitting full $p \times p$ covariance matrices for each mixture component when there are $n_g < p$ observations for the component. One purely numerical approach we have used to mitigate singularities was to regularize the updated covariance matrices at each iteration in the EM algorithm using a suitably chosen eigenvalue cutoff. Future work will include using analogues of the GPCM models, via eigen-decomposition of the component scale matrices as well as investigating alternative approaches for dealing with singularity problems. The latter might include working on distinct subspaces for fitting each covariance matrix and only injecting into a common subspace when forming \mathbf{M} .

The real data sets used for our illustrations were selected because they were previously used to illustrate the performance of some of the comparator methods. Therefore, it is encouraging that HMMDR consistently outperformed the comparator methods on real data sets. As part of the development of a companion R package, we will study whether multiple components should be available to represent a class in HMMDR discriminant analysis. Finally, the application of our approach within the fractionally-supervised classification framework (Vrbik and McNicholas, 2013) will also be a subject of future work.

Acknowledgements

The authors are grateful to an associate editor and two anonymous reviewers for their very helpful comments. This work was supported by an Ontario Graduate Scholarship (Morris), an Early Researcher Award from the Government of Ontario (McNicholas), and a Discovery Grant from the Natural Sciences and Engineering Research Council of Canada (NSERC; McNicholas). The computing equipment used was provided through a Research Tools and Instruments Grant from NSERC.

References

- Abramowitz, M., Stegun, I., 1972. Handbook of Mathematical Functions, 9th Edition. Dover, New York.
- Airoldi, J.-P., Hoffmann, R. S., 1984. Age variation in voles (*Microtus californicus*, *M. ochrogaster*) and its significance for systematic studies. Occasional papers of the Museum of Natural History, University of Kansas, Lawrence KS 111, 1–45.
- Aitken, A. C., 1926. On Bernoulli’s numerical solution of algebraic equations. Proceedings of the Royal Society of Edinburgh 46, 289–305.
- Alon, U., Barkai, N., Notterman, D., Gish, K., Ybarra, S., Mack, D., Levine, A., 1999. Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. Proceedings of the National Academy of Sciences 96 (12), 6745–6750.
- Andrews, J. L., McNicholas, P., 2012a. Model-based clustering, classification, and discriminant analysis via mixtures of multivariate *t*-distributions: The *t*EIGEN family. Statistics and Computing 22 (5), 1021–1029.
- Andrews, J. L., McNicholas, P. D., 2011. Mixtures of modified *t*-factor analyzers for model-based clustering, classification, and discriminant analysis. Journal of Statistical Planning and Inference 141 (4), 1479–1486.

- Andrews, J. L., McNicholas, P. D., 2012b. teigen: Model-based clustering and classification with the multivariate t -distribution. R package version 1.0.
- Andrews, J. L., McNicholas, P. D., Subedi, S., 2011. Model-based classification via mixtures of multivariate t -distributions. *Computational Statistics and Data Analysis* 55 (1), 520–529.
- Baek, J., McLachlan, G. J., Flack, L. K., 2009. mcfa: Fits mixtures of common factor analyzers to a given data set. R package version 1.0.2.
- Baek, J., McLachlan, G. J., Flack, L. K., 2010. Mixtures of factor analyzers with common factor loadings: applications to the clustering and visualization of high-dimensional data. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 32 (7), 1298–1309.
- Banfield, J. D., Raftery, A. E., 1993. Model-based Gaussian and non-Gaussian clustering. *Biometrics* 49 (3), 803–821.
- Barndorff-Nielsen, O., 1977. Exponentially decreasing distributions for the logarithm of particle size. *Proceedings of the Royal Society A* 353, 401–419.
- Barndorff-Nielsen, O., Halgreen, C., 1977. Infinite divisibility of the hyperbolic and generalized inverse Gaussian distributions. *Z. Wahrscheinlichkeitstheorie verw. Gebiete* 38, 309–311.
- Baum, L. E., Petrie, T., Soules, G., Weiss, N., 1970. A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains. *Annals of Mathematical Statistics* 41, 164–171.
- Bensmail, H., Celeux, G., 1996. Regularized Gaussian discriminant analysis through eigenvalue decomposition. *Journal of the American Statistical Association* 91, 1743–1748.
- Blæsild, P., 1978. The shape of the generalized inverse Gaussian and hyperbolic distributions. Research Report 37, Aarhus University, Denmark, Department of Theoretical Statistics.
- Böhning, D., Dietz, E., Schaub, R., Schlattmann, P., Lindsay, B., 1994. The distribution of the likelihood ratio for mixtures of densities from the one-parameter exponential family. *Annals of the Institute of Statistical Mathematics* 46, 373–388.
- Boulesteix, A.-L., Lambert-Lacroix, S., Peyre, J., Strimmer, K., 2011. plsgenomics: PLS analyses for genomics. R package version 1.2-6.
- Bouveyron, C., Brunet, C., 2012. Simultaneous model-based clustering and visualization in the Fisher discriminative subspace. *Statistics and Computing* 22 (1), 301–324.
- Browne, R. P., McNicholas, P. D., 2015. A mixture of generalized hyperbolic distributions. *Canadian Journal of Statistics* 43 (2), 176–198.
- Celeux, G., Govaert, G., 1995. Gaussian parsimonious clustering models. *Pattern Recognition* 28, 781–793.
- Dean, N., Murphy, T. B., Downey, G., 2006. Using unlabelled data to update classification rules with applications in food authenticity studies. *Journal of the Royal Statistical Society: Series C* 55 (1), 1–14.
- Dempster, A. P., Laird, N. M., Rubin, D. B., 1977. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society: Series B* 39 (1), 1–38.
- Everitt, B. S., Hand, D. J., 1981. *Finite Mixture Distributions*. Chapman and Hall, London.
- Faraway, J., 2011. faraway: Functions and datasets for books by Julian Faraway. R package version 1.0.5.
- Flury, B., 2010. Flury: Data Sets from Flury, 1997. R package version 0.1-3.
- Flury, B., Riedwyl, H., 1988. *Multivariate Statistics: A Practical Approach*. Cambridge University Press.

- Flury, B. D., 1997. *A First Course in Multivariate Statistics*. Springer, New York.
- Forina, M., Armanino, C., Castino, M., Ubigli, M., 1986. Multivariate data analysis as a discriminating method of the origin of wines. *Vitis* 25, 189–201.
- Fraley, C., Raftery, A. E., 1999. MCLUST: Software for model-based cluster analysis. *Journal of Classification* 16, 297–306.
- Fraley, C., Raftery, A. E., 2002. Model-based clustering, discriminant analysis, and density estimation. *Journal of the American Statistical Association* 97 (458), 611–631.
- Fraley, C., Raftery, A. E., Murphy, T. B., Scrucca, L., 2012. MCLUST version 4 for R: Normal mixture modeling for model-based clustering, classification, and density estimation. Department of Statistics, University of Washington.
- Franczak, B. C., Browne, R. P., McNicholas, P. D., 2014. Mixtures of shifted asymmetric Laplace distributions. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 36 (6), 1149–1157.
- Franczak, B. C., Tortora, C., Browne, R. P., McNicholas, P. D., 2015. Unsupervised learning via mixtures of skewed distributions with hypercube contours. *Pattern Recognition Letters* 58 (1), 69–76.
- Greselin, F., Ingrassia, S., 2010. Constrained monotone EM algorithms for mixtures of multivariate t -distributions. *Statistics and Computing* 20 (1), 9–22.
- Härdle, W. K., Simar, L., 2011. *Applied Multivariate Statistical Analysis*. Springer.
- Hastie, T., Tibshirani, R., 1996. Discriminant analysis by Gaussian mixtures. *Journal of the Royal Statistical Society: Series B* 58 (1), 155–176.
- Hubert, L., Arabie, P., 1985. Comparing partitions. *Journal of Classification* 2, 193–218.
- Hubert, M., Rousseeuw, P. J., Vanden Branden, K., 2005. ROBPCA: a new approach to robust principal components analysis. *Technometrics* 47, 64–79.
- Hurley, C., 2004. Clustering visualizations of multivariate data. *Journal of Computational and Graphical Statistics* 13 (4), 788–806.
- Hurley, C., 2010. *gclus: Clustering Graphics*. R package version 1.3.
- Lee, S. X., McLachlan, G. J., 2013. On mixtures of skew normal and skew t -distributions. [arXiv:1211.3602v2](https://arxiv.org/abs/1211.3602v2).
- Li, K. C., 1991. Sliced inverse regression for dimension reduction (with discussion). *Journal of the American Statistical Association* 86, 316–342.
- Li, K. C., 2000. High dimensional data analysis via the SIR/PHD approach, unpublished manuscript.
- Lin, T.-I., 2010. Robust mixture modeling using multivariate skew t -distributions. *Statistics and Computing* 20, 343–356.
- Lindsay, B. G., 1995. Mixture models: Theory, geometry and applications. In: *NSF-CBMS Regional Conference Series in Probability and Statistics*. Vol. 5. Hayward, California: Institute of Mathematical Statistics.
- Mangasarian, O. L., Street, W. N., Wolberg, W. H., 1995. Breast cancer diagnosis and prognosis via linear programming. *Operations Research* 43 (4), 570–577.
- McLachlan, G. J., 1982. The classification and mixture maximum likelihood approaches to cluster analysis. Vol. 2 of *Handbook of Statistics*. North-Holland, Amsterdam, pp. 199–208.
- McLachlan, G. J., Basford, K. E., 1988. *Mixture Models: Inference and applications to clustering*. Marcel Dekker Inc., New York.

- McLachlan, G. J., Bean, R. W., Peel, D., 2002. Mixture model-based approach to the clustering of microarray expression data. *Bioinformatics* 18, 413–422.
- McLachlan, G. J., Peel, D., 1998. Robust cluster analysis via mixtures of multivariate t -distributions. In: *Lecture Notes in Computer Science*. Vol. 1451. Springer-Verlag, Berlin, pp. 658–666.
- McLachlan, G. J., Peel, D., 2000. *Finite Mixture Models*. John Wiley & Sons, New York.
- McNeil, A. J., Frey, R., Embrechts, P., 2005. *Quantitative Risk Management: Concepts, Techniques and Tools*. Princeton University Press.
- McNicholas, P. D., 2010. Model-based classification using latent Gaussian mixture models. *Journal of Statistical Planning and Inference* 140, 1175–1181.
- McNicholas, P. D., 2013. Model-based clustering and classification via mixtures of multivariate t -distributions. In: Guidici, P., Ingrassia, S., Vichi, M. (Eds.), *Statistical Models for Data Analysis. Studies in Classification, Data Analysis, and Knowledge Organization*. Springer International Publishing Switzerland, pp. 233–240.
- McNicholas, P. D., Jampani, K. R., McDaid, A. F., Murphy, T. B., Banks, L., 2011. *pgmm: Parsimonious Gaussian Mixture Models*. R package version 1.0.
- McNicholas, P. D., Murphy, T. B., 2008. Parsimonious Gaussian mixture models. *Statistics and Computing* 18, 285–296.
- McNicholas, P. D., Murphy, T. B., 2010. Model-based clustering of microarray expression data via latent Gaussian mixture models. *Bioinformatics* 26 (21), 2705–2712.
- McNicholas, P. D., Murphy, T. B., McDaid, A. F., Frost, D., 2010. Serial and parallel implementations of model-based clustering via parsimonious Gaussian mixture models. *Computational Statistics and Data Analysis* 54 (3), 711–723.
- McNicholas, S. M., McNicholas, P. D., Browne, R. P., 2013. Mixtures of variance-gamma distributions. Arxiv preprint arXiv:1309.2695.
- Morris, K., McNicholas, P. D., 2013. Dimension reduction for model-based clustering via mixtures of shifted asymmetric Laplace distributions. *Statistics and Probability Letters* 83, 2088–2093.
- Morris, K., McNicholas, P. D., Scrucca, L., 2013. Dimension reduction for model-based clustering via mixtures of multivariate t -distributions. *Advances in Data Analysis and Classification*, 1–18.
- Murray, P. M., Browne, R. B., McNicholas, P. D., 2014a. Mixtures of skew- t factor analyzers. *Computational Statistics and Data Analysis* 77, 326–335.
- Murray, P. M., McNicholas, P. D., Browne, R. B., 2014b. A mixture of common skew- t factor analyzers. *Stat* 3 (1), 68–82.
- Orchard, T., Woodbury, M. A., 1972. A missing information principle: theory and applications. In: Le Cam, L. M., Neyman, J., Scott, E. L. (Eds.), *Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Theory of Statistics*. University of California Press, Berkeley, pp. 697–715.
- Peel, D., McLachlan, G. J., 2000. Robust mixture modelling using the t -distribution. *Statistics & Computing* 10, 339–348.
- Qiu, W.-L., Joe, H., 2006. Generation of random clusters with specified degree of separation. *Journal of Classification* 23 (2), 315–334.
- R Core Team, 2013. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
URL <http://www.R-project.org>

- Raftery, A. E., Dean, N., 2006. Variable selection for model-based clustering. *Journal of the American Statistical Association* 101 (473), 168–178.
- Rand, W. M., 1971. Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical Association* 66, 846–850.
- Schwarz, G., 1978. Estimating the dimension of a model. *Annals of Statistics* 6, 461–464.
- Schwender, H., 2012. siggenes: Multiple testing using SAM and Efron’s empirical Bayes approaches. R package version 1.32.0.
- Scrucca, L., 2010. Dimension reduction for model-based clustering. *Statistics & Computing* 20 (4), 471–484.
- Scrucca, L., 2013. Graphical tools for model-based mixture discriminant analysis. *Advances in Data Analysis and Classification*, 1–19.
- Scrucca, L., Raftery, A. E., Dean, N., 2013. clustvarsel: A package implementing variable selection for model-based clustering in R, submitted to *Journal of Statistical Software*.
- Steane, M. A., McNicholas, P. D., Yada, R., 2012. Model-based classification via mixtures of multivariate t-factor analyzers. *Communications in Statistics – Simulation and Computation* 41 (4), 510–523.
- Sundberg, R., 1974. Maximum likelihood theory for incomplete data from an exponential family. *Scandinavian Journal of Statistics* 1, 49–58.
- Todorov, V., Filzmoser, P., 2009. An object-oriented framework for robust multivariate analysis. *Journal of Statistical Software* 32 (3), 1–47.
- Tortora, C., Franczak, B. C., Browne, R. P., McNicholas, P. D., 2015. A mixture of coalesced generalized hyperbolic distributions. arXiv preprint arXiv:1403.2332v6.
- Tusher, V., Tibshirani, R., Chu, G., 2001. Significance analysis of microarrays applied to the ionizing radiation response. *Proceedings of the National Academy of Sciences* 98 (9), 5116–5121.
- Vrbik, I., McNicholas, P. D., 2012. Analytic calculations for the EM algorithm for multivariate skew-mixture models. *Statistics and Probability Letters* 82 (6), 1169–1174.
- Vrbik, I., McNicholas, P. D., 2013. Fractionally-supervised classification. arXiv:1307.3598v2.
- Vrbik, I., McNicholas, P. D., 2014. Parsimonious skew mixture models for model-based clustering and classification. *Computational Statistics and Data Analysis* 71, 196–210.
- Wolfe, J. H., 1963. Object cluster analysis of social areas. Master’s thesis, University of California, Berkeley.

A Selected Genes for the Colon Cancer Data

Table 19: Genes selected with SAM for the colon data.

Hsa.462 (Human serine kinase mRNA)
Hsa.549 (Transcription factor IIIA)
Hsa.601 (Human aspartyl-tRNA synthetase alpha-2 subunit mRNA)
Hsa.627 (Human monocyte-derived neutrophil-activating protein (MONAP) mRNA)
Hsa.773 (Macrophage migration inhibitory factor (Human))
Hsa.821 (Human hmgI mRNA for high mobility group protein Y)
Hsa.831 (Mitochondrial matrix protein P1 precursor (Human))
Hsa.957 (Human nucleolar protein (B23) mRNA)
Hsa.1832 (Myosin regulatory light chain 2, smooth muscle isoform (Human))
Hsa.2097 (Human vasoactive intestinal peptide (VIP) mRNA)
Hsa.2645 (H.sapiens ckshs2 mRNA for Cks1 protein homologue)
Hsa.2928 (H.sapiens mRNA for p cadherin)
Hsa.3016 (S-100P protein (Human))
Hsa.3306 (Human gene for heterogeneous nuclear ribonucleoprotein (hnRNP) core protein A1)
Hsa.3331 (Nucleoside diphosphate kinase A (Human))
Hsa.5971 (Human splicing factor SRp30c mRNA)
Hsa.6472 (Tubulin beta chain (Haliotis discus))
Hsa.6814 (Collagen alpha 2(XI) chain (Homo sapiens))
Hsa.8125 (Human)
Hsa.8147 (Human desmin gene)
Hsa.36689 (H.sapiens mRNA for GCAP-II/uroguanylin precursor)
Hsa.36952 (Complement factor D precursor (Homo sapiens))
Hsa.37937 (Myosyn heavy chain, nonmuscle (Gallus gallus))
