



## Capabilities of R Package `mixAK` for Clustering Based on Multivariate Continuous and Discrete Longitudinal Data

Arnošt Komárek  
Charles University in Prague

Lenka Komárková  
University of Economics in Prague

---

### Abstract

R package `mixAK` originally implemented routines primarily for Bayesian estimation of finite normal mixture models for possibly interval-censored data. The functionality of the package was considerably enhanced by implementing methods for Bayesian estimation of mixtures of multivariate generalized linear mixed models proposed in Komárek and Komárková (2013). Among other things, this allows for a cluster analysis (classification) based on multivariate continuous and discrete longitudinal data that arise whenever multiple outcomes of a different nature are recorded in a longitudinal study. This package also allows for a data-driven selection of a number of clusters as methods for selecting a number of mixture components were implemented. A model and clustering methodology for multivariate continuous and discrete longitudinal data is overviewed. Further, a step-by-step cluster analysis based jointly on three longitudinal variables of different types (continuous, count, dichotomous) is given, which provides a user manual for using the package for similar problems.

*Keywords:* cluster analysis, generalized linear mixed model, functional data, multivariate longitudinal data, R package.

---

## 1. Introduction

It is a common practice in longitudinal studies to gather multiple outcomes, both continuous and discrete at each subject's follow-up visit leading to multivariate longitudinal data. In many research areas, the interest then lies in classifying (clustering) the subjects into groups (clusters) on the basis of such multivariate longitudinal data. To be more specific, we first introduce a basic notation. Suppose that  $\mathbf{t}_i = (t_{i,1}, \dots, t_{i,n_i})^\top$ ,  $i = 1, \dots, N$ , are the visit times of  $N$  subjects involved in a longitudinal study, where  $n_i$  is the  $i$ th subject's number

of visits which may differ between subjects, as well as the actual visit times. Further, let  $\mathbf{Y}_{i,r} = (Y_{i,r,1}, \dots, Y_{i,r,n_i})^\top = (Y_{i,r,1}(t_{i,1}), \dots, Y_{i,r,n_i}(t_{i,n_i}))^\top$ ,  $r = 1, \dots, R$ , be the random vectors representing the  $i$ th subject's longitudinal measurements of the  $r$ th outcome being recorded. The problem of clustering based on multivariate longitudinal data lies in using complete observational vectors  $\mathbf{Y}_i = (\mathbf{Y}_{i,1}^\top, \dots, \mathbf{Y}_{i,R}^\top)^\top$  together with  $\mathbf{t}_i$ ,  $i = 1, \dots, N$ , and possibly other exogenous covariates for classification of subjects into one of the  $K$  groups where  $K$  is either known or unknown.

### 1.1. Model-based clustering

Model-based clustering became a popular method of classification in situations where it is suitable to distinguish the  $K$  clusters by different probabilistic models (Bock 1996; Fraley and Raftery 2002). Initially, we assume that  $K$  is known. As usual in this context, we introduce the unobservable component allocations  $U_1, \dots, U_N \in \{1, \dots, K\}$ ,

$$P(U_i = k; \mathbf{w}) = w_k, \quad i = 1, \dots, N, \quad k = 1, \dots, K,$$

where  $\mathbf{w} = (w_1, \dots, w_K)^\top$  is a vector of unknown cluster proportions which are positive and sum to unity. The meaning of the component allocations is so that  $U_i = k$  when the  $i$ th subject's observational random vector  $\mathbf{Y}_i$  was generated by the  $k$ th probabilistic model represented by a model density  $f_{i,k}(\mathbf{y}_i; \boldsymbol{\xi}_k, \boldsymbol{\xi})$ , where  $\boldsymbol{\xi}_k$  is a vector of cluster-specific and  $\boldsymbol{\xi}$  a vector of common unknown parameters. The subscript  $i$  in  $f_{i,k}$ , which is a conditional density of  $\mathbf{Y}_i$  given  $U_i = k$ , points to the fact that  $f_{i,k}$  may depend on subject specific factors like the visit times  $\mathbf{t}_i$  or other covariates, and allows us to consider also regression models as cluster characteristics. The marginal density of  $\mathbf{Y}_i$ , and hence also the likelihood contribution of the  $i$ th subject, is then a mixture density

$$f_i(\mathbf{y}_i; \boldsymbol{\theta}) = \sum_{k=1}^K w_k f_{i,k}(\mathbf{y}_i; \boldsymbol{\xi}_k, \boldsymbol{\xi}), \quad (1)$$

where  $\boldsymbol{\theta} = (\mathbf{w}^\top, \boldsymbol{\xi}_1^\top, \dots, \boldsymbol{\xi}_K^\top, \boldsymbol{\xi}^\top)^\top$  denotes a vector of all unknown model parameters. Model-based clustering is then based on the estimated values  $\hat{p}_{i,k}$ ,  $i = 1, \dots, N$ ,  $k = 1, \dots, K$ , of the individual component probabilities

$$p_{i,k} = p_{i,k}(\boldsymbol{\theta}) = P(U_i = k \mid \mathbf{Y}_i = \mathbf{y}_i; \boldsymbol{\theta}) = \frac{w_k f_{i,k}(\mathbf{y}_i; \boldsymbol{\xi}_k, \boldsymbol{\xi})}{f_i(\mathbf{y}_i; \boldsymbol{\theta})}, \quad i = 1, \dots, N, \quad k = 1, \dots, K. \quad (2)$$

Different classification rules are possible depending on the specific choice of loss function associated with a misclassification. One possible classification rule is to assign subject  $i$  into the cluster  $g(i)$  such that  $\hat{p}_{i,g(i)} = \max_{k=1, \dots, K} \hat{p}_{i,k}$ .

### 1.2. R package mixAK

The R (R Core Team 2014) package **mixAK** (Komárek 2014), which is available from the Comprehensive R Archive Network (CRAN) at <http://CRAN.R-project.org/package=mixAK>, started as a set of routines for Bayesian estimation of finite mixture models and subsequent

clustering based on possibly censored data. In the earlier version of the package, which is described in Komárek (2009), only the simple case was covered where the densities  $f_{i,k}$  in (1) were all assumed to be multivariate normal with means and covariance matrices depending on  $k$  only, i.e.,  $f_{i,k} = f_k \equiv \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ ,  $\boldsymbol{\xi}_k = (\boldsymbol{\mu}_k^\top, \text{vec}^\top(\boldsymbol{\Sigma}_k))^\top$ ,  $k = 1, \dots, K$ .

Following two methodological papers (Komárek, Hansen, Kuiper, van Buuren, and Lesaffre 2010; Komárek and Komárková 2013), the package has been considerably extended as follows. It is assumed that the densities  $f_{i,k}$  correspond to multivariate generalized linear mixed models, which leads to the density  $f_i$  (Equation 1) being a mixture of multivariate generalized linear mixed models. We describe this concept in more detail in Section 2. These extensions result in two new principal capabilities of the package that are only rarely implemented elsewhere. These are:

- (i) joint analysis of multivariate longitudinal data, both continuous and discrete, by the mean of an extension of a well understood generalized linear mixed model;
- (ii) cluster analysis based on multivariate longitudinal data.

Even though we shall exemplify the whole methodology on problems from the area of longitudinal studies, all methods as well as the **mixAK** package have much broader area of applications in fact covering data exhibiting all types of repeated measurements. For example, they can be directly applied to perform cluster analysis for functional data, to name one.

In this manuscript, we concentrate on introducing the new capabilities of the package **mixAK** related to perform cluster analysis for multivariate longitudinal data in particular. To this end, we first finalize the introduction by a brief overview of other possible R implementations of methods which might be considered for clustering based on longitudinal data. Nevertheless, we also explain that most of these apparent competitors to our package can only be used for clustering longitudinal data under restrictive assumptions not requested by our package. Further, in Section 2, we first overview the particular form of the mixture model (1) that underlies the methods implemented in the package **mixAK** and then discuss how to use the model for clustering and also how to infer a suitable number of mixture components. The main part of the paper is given in Section 3, where we provide a detailed R example. It is a cluster analysis based on a medical dataset involving three longitudinal outcomes of different nature (continuous, count, dichotomous), which served as a motivating example for Komárek and Komárková (2013). This paper finishes with conclusions and an outlook in Section 4.

### 1.3. Possible competitors

A variety of model-based clustering methods have been implemented and the methods are available as contributed R packages. If the longitudinal data at hand are regularly sampled ( $n_1 = \dots = n_N = n$ ,  $\mathbf{t}_1 = \dots = \mathbf{t}_N = \mathbf{t}$ ), it is reasonable to assume that the vectors  $\mathbf{Y}_1, \dots, \mathbf{Y}_N$  are not only independent but also identically distributed, i.e., i.i.d. Consequently, it is possible to consider the clustering methods based on finite mixtures of classical distributions. Besides the earlier version of the **mixAK** package, these are available in the R packages **mclust** (Fraley, Raftery, Murphy, and Scrucca 2012) or **teigen** (Andrews and McNicholas 2012, 2013), for instance, where the mixture components are assumed to follow either a multivariate normal or t-distributions. Mixtures with not only continuous but also discrete components

are implemented in the R packages **mixdist** (Macdonald and Du 2012) or **Rmixmod** (Auder, Le Bret, Iovleff, and Langrognet 2014). Cluster analysis with high dimensional data can also be performed by the means of the R package **HDclassif** (Bergé, Bouveyron, and Girard 2012). Nevertheless, for all of the above-mentioned approaches, the mixture parameters are not allowed in a structural way to depend on exogenous factors like the visit time vectors  $\mathbf{t}_1, \dots, \mathbf{t}_N$ , for example. This makes them impractical or inefficient for actual longitudinal data. The possibility to model the dependence of mixture parameters on factors like the visit times  $\mathbf{t}_1, \dots, \mathbf{t}_N$  and use the resulting model for clustering is provided by the R packages **long-clust** (McNicholas, Jampani, and Subedi 2012) or **MFDA** (Zhong and Ma 2012). However, only regularly sampled, and only continuous longitudinal data can be clustered using those packages.

In biostatistical applications, as well as in other research areas, the longitudinal data are typically irregularly sampled, i.e., having in general different values of numbers of visits  $n_1, \dots, n_N$  and/or different visit time vectors  $\mathbf{t}_1, \dots, \mathbf{t}_N$ . Model-based clustering methods for such data suggested in the literature are then usually based on a mixture of suitable regression models. A mixture of linear mixed models is the basis for clustering of continuous longitudinal data, the approach being implemented in the R package **mixtools** (Benaglia, Chauveau, Hunter, and Young 2009). Both continuous and discrete longitudinal data can be clustered via mixtures of various types of mixed models using the R package **lcmm** (Proust-Lima, Philipps, Diakite, and Liqueur 2014). However, it is not possible to use *jointly* continuous and discrete outcomes in one analysis. This rules out this package for clustering based on general multivariate longitudinal data, where different distributional assumptions, e.g., Gaussian for elements of  $\mathbf{Y}_{i,1}$ , and Bernoulli for elements of  $\mathbf{Y}_{i,2}$ ,  $i = 1, \dots, N$ , might be unavoidable. According to the best of our knowledge, the only R package allowing for cluster analysis based on actual multivariate longitudinal data of different nature (both continuous and discrete) is **flexmix** (Grün and Leisch 2008). It implements a mixture of generalized linear models, possibly with repeated measurements estimated using likelihood principles by the means of the EM algorithm. Nevertheless, in the case of multivariate responses (i.e.,  $R > 1$  in our notation), the model used by **flexmix** assumes that these are independent for  $r = 1, \dots, R$  which might be unrealistic.

The mixture of multivariate generalized linear mixed models provides some features, which make the **mixAK** package different and to some extent more broadly applicable compared to the above-mentioned implementations. Briefly, those features are such that they allow for (i) irregularly sampled longitudinal data; (ii) multivariate longitudinal data with responses of different nature (continuous and discrete); and (iii) multivariate responses are not necessarily assumed to be independent for one subject.

## 2. Model and clustering procedure

Methodology which underlies the procedures for clustering based on multivariate longitudinal data implemented in the package **mixAK** is described in detail in Komárek and Komárková (2013) and its electronic supplement. Nevertheless, to make this paper relatively standalone and to introduce the necessary notation, we provide a brief overview in this section.

### 2.1. Multivariate mixture generalized linear mixed model

The mixture model (1) assumed by the package **mixAK** for the observable random vectors

$\mathbf{Y}_i = (Y_{i,1,1}, \dots, Y_{i,R,n_i})^\top$ ,  $i = 1, \dots, N$ , as a basis for the clustering procedure is determined by the following assumptions.

- (A1) For each  $i = 1, \dots, N$ ,  $r = 1, \dots, R$  and each  $j = 1, \dots, n_i$ ,  $Y_{i,r,j}$  follows a distribution  $\mathcal{D}_r$  from the exponential family with the dispersion parameter  $\phi_r$  (fixed or unknown depending on the considered exponential family) and the mean given by

$$h_r^{-1}\{\mathbb{E}(Y_{i,r,j} \mid \mathbf{B}_{i,r} = \mathbf{b}_{i,r}; \boldsymbol{\alpha}_r)\} = \mathbf{x}_{i,r,j}^\top \boldsymbol{\alpha}_r + \mathbf{z}_{i,r,j}^\top \mathbf{b}_{i,r}, \quad (3)$$

where

- (i)  $h_r^{-1}$  is a chosen link function;
- (ii)  $\mathbf{x}_{i,r,j}^\top \in \mathbb{R}^{p_r}$ ,  $\mathbf{z}_{i,r,j}^\top \in \mathbb{R}^{q_r}$  are vectors of known covariates (visit times  $t_{i,j}$  and possibly other factors) such that the matrix  $(\mathbb{X}_r \mathbb{Z}_r)$ , where

$$\mathbb{X}_r = \begin{pmatrix} \mathbf{x}_{1,r,1}^\top \\ \vdots \\ \mathbf{x}_{N,r,n_N}^\top \end{pmatrix}, \quad \mathbb{Z}_r = \begin{pmatrix} \mathbf{z}_{1,r,1}^\top \\ \vdots \\ \mathbf{z}_{N,r,n_N}^\top \end{pmatrix},$$

is of full column rank  $p_r + q_r$ ;

- (iii)  $\boldsymbol{\alpha}_r \in \mathbb{R}^{p_r}$  is a vector of unknown parameters (fixed effects);
  - (iv)  $\mathbf{B}_{i,r} \in \mathbb{R}^{q_r}$  is a random vector (random effects).
- (A2) For each  $i = 1, \dots, N$ , the conditional distribution of the joint (over  $R$  markers) random effects vector  $\mathbf{B}_i = (\mathbf{B}_{i,1}^\top, \dots, \mathbf{B}_{i,R}^\top)^\top \in \mathbb{R}^q$ ,  $q = \sum_{r=1}^R q_r$ , given  $U_i = k$  (given the  $i$ th subject belongs to the  $k$ th group) is a (multivariate) normal with unknown mean  $\boldsymbol{\mu}_k \in \mathbb{R}^q$  and unknown  $q \times q$  positive definite covariance matrix  $\mathbb{D}_k$ ,  $k = 1, \dots, K$ , i.e.,

$$\mathbf{B}_i \mid U_i = k \sim \mathcal{N}_q(\boldsymbol{\mu}_k, \mathbb{D}_k), \quad k = 1, \dots, K. \quad (4)$$

- (A3) For each  $i = 1, \dots, N$ , the random variables  $Y_{i,1,1}, \dots, Y_{i,R,n_i}$  are conditionally independent given the random effects vector  $\mathbf{B}_i$ .

- (A4) Random vectors  $\mathbf{Y}_1, \dots, \mathbf{Y}_N$  are independent.

- (A5) Random effects vectors  $\mathbf{B}_1, \dots, \mathbf{B}_N$  are independent.

In summary, the cluster-specific model parameters  $\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_K$  are composed of the means and covariance matrices of the conditional distributions of random effects, i.e.,

$$\boldsymbol{\xi}_k = (\boldsymbol{\mu}_k^\top, \text{vec}^\top(\mathbb{D}_k))^\top, \quad k = 1, \dots, K.$$

The model parameters common to all clusters are the fixed effects and the dispersion parameters, i.e.,

$$\boldsymbol{\xi} = (\boldsymbol{\alpha}_1^\top, \dots, \boldsymbol{\alpha}_R^\top, \phi_1, \dots, \phi_R)^\top. \quad (5)$$

The cluster specific model density  $f_{i,k}$ ,  $i = 1, \dots, N$ ,  $k = 1, \dots, K$ , is then

$$f_{i,k}(\mathbf{y}_i; \boldsymbol{\xi}_k, \boldsymbol{\xi}) = \int_{\mathbb{R}^q} \left\{ \prod_{r=1}^R \prod_{j=1}^{n_i} f_{\mathcal{D}_r}(y_{i,r,j}; \boldsymbol{\alpha}_r, \phi_r, \mathbf{b}_{i,r}) \right\} \varphi(\mathbf{b}_i; \boldsymbol{\mu}_k, \mathbb{D}_k) d\mathbf{b}_i, \quad (6)$$

where  $f_{\mathcal{D}_r}$  is the exponential family density following from the assumption (A1). Further,  $\varphi(\cdot; \boldsymbol{\mu}_k, \mathbb{D}_k)$  is a density of the (multivariate) normal distribution with a mean  $\boldsymbol{\mu}_k$  and a covariance matrix  $\mathbb{D}_k$ ,  $k = 1, \dots, K$ , following from assumption (A2).

With  $R = 1$ , the model density  $f_{i,k}$  in (6) is the  $i$ th subject's likelihood contribution as if the generalized linear mixed model (GLMM) with  $\mathcal{N}(\boldsymbol{\mu}_k, \mathbb{D}_k)$  distributed random effects is assumed for the observed data. With  $R > 1$ , a multivariate generalized linear mixed model (MGLMM) is obtained where the dependence among the random vectors  $\mathbf{Y}_{i,1}, \dots, \mathbf{Y}_{i,R}$  representing different markers is induced by non-diagonal covariance matrix  $\mathbb{D}_k$  of the random effects vector  $\mathbf{B}_i$  in general. Finally, substituting  $f_{i,k}$  from (6) into (1) lead to the  $i$ th subject's likelihood contribution

$$f_i(\mathbf{y}_i; \boldsymbol{\theta}) = \sum_{k=1}^K w_k \int_{\mathbb{R}^q} \left\{ \prod_{r=1}^R \prod_{j=1}^{n_i} f_{\mathcal{D}_r}(y_{i,r,j}; \boldsymbol{\alpha}_r, \phi_r, \mathbf{b}_{i,r}) \right\} \varphi(\mathbf{b}_i; \boldsymbol{\mu}_k, \mathbb{D}_k) d\mathbf{b}_i, \quad (7)$$

$$= \int_{\mathbb{R}^q} \left\{ \prod_{r=1}^R \prod_{j=1}^{n_i} f_{\mathcal{D}_r}(y_{i,r,j}; \boldsymbol{\alpha}_r, \phi_r, \mathbf{b}_{i,r}) \right\} \left\{ \sum_{k=1}^K w_k \varphi(\mathbf{b}_i; \boldsymbol{\mu}_k, \mathbb{D}_k) \right\} d\mathbf{b}_i. \quad (8)$$

It follows from the expression above that the model assumed for the observable random vectors  $\mathbf{Y}_1, \dots, \mathbf{Y}_N$  can be interpreted either as a mixture of multivariate generalized linear mixed models with normally distributed random effects (Equation 7), or as a multivariate generalized linear mixed model with a normal mixture in the random effects distribution (Equation 8), where the overall mean and the overall covariance matrix of the random effects  $\mathbf{B}_i$ ,  $i = 1, \dots, N$ , are given by

$$\boldsymbol{\beta} = \mathbb{E}(\mathbf{B}_i; \boldsymbol{\theta}) = \sum_{k=1}^K w_k \boldsymbol{\mu}_k, \quad (9)$$

$$\mathbb{D} = \text{VAR}(\mathbf{B}_i; \boldsymbol{\theta}) = \sum_{k=1}^K w_k \left\{ \mathbb{D}_k + \left( \boldsymbol{\mu}_k - \sum_{j=1}^K w_j \boldsymbol{\mu}_j \right) \left( \boldsymbol{\mu}_k - \sum_{j=1}^K w_j \boldsymbol{\mu}_j \right)^\top \right\}. \quad (10)$$

Consequently, we call our model a multivariate mixture generalized linear mixed model (MMGLMM). Following the above considerations, the model likelihood and its observed data deviance are

$$L(\boldsymbol{\theta}) = \prod_{i=1}^N f_i(\mathbf{y}_i; \boldsymbol{\theta}), \quad D(\boldsymbol{\theta}) = -2 \log \{L(\boldsymbol{\theta})\}, \quad (11)$$

respectively.

Finally, we point out that in the package **mixAK**, the following exponential distributions  $\mathcal{D}_r$  and the link functions  $h_r^{-1}$  from assumption (A1) are implemented: (a) Gaussian with the identity link, i.e., a linear mixed model for the  $r$ th marker where the dispersion parameter  $\phi_r$  is the unknown residual variance; (b) Poisson with the log link where  $\phi_r = 1$ ; (c) Bernoulli with the logit link where again  $\phi_r = 1$ .

## 2.2. Bayesian inference

For largely computational reasons, the Bayesian approach based on the output from the Markov chain Monte Carlo (MCMC) simulation is exploited to infer the unknown model

parameter vector  $\boldsymbol{\theta}$  and to perform the clustering. In the sequel, let  $p(\cdot)$  and  $p(\cdot|\cdot)$  be generic symbols for (conditional) distributions and densities. As is usual in Bayesian statistics, the latent quantities (random effects  $\mathbf{B}_i$  and component allocations  $U_i$ ,  $i = 1, \dots, N$ ) are considered as additional model parameters with the joint prior distribution for all model parameters specified hierarchically following the structure of the model outlined in Sections 1.1 and 2.1 as

$$\begin{aligned} p(\boldsymbol{\theta}, \mathbf{b}_1, \dots, \mathbf{b}_N, u_1, \dots, u_N) &= \prod_{i=1}^N \left\{ p(\mathbf{b}_i | u_i, \boldsymbol{\theta}) p(u_i | \boldsymbol{\theta}) \right\} p(\boldsymbol{\theta}) \\ &= \prod_{i=1}^N \left\{ \varphi(\mathbf{b}_i; \boldsymbol{\mu}_{u_i}, \mathbb{D}_{u_i}) w_{u_i} \right\} p(\boldsymbol{\theta}). \end{aligned}$$

The prior distribution  $p(\boldsymbol{\theta})$  of the primary model parameters is then specified to be weakly informative. Finally, the MCMC methods are used to generate a sample

$$\mathcal{S}_M = \left\{ \left( \boldsymbol{\theta}^{(m)}, \mathbf{b}_1^{(m)}, \dots, \mathbf{b}_N^{(m)}, u_1^{(m)}, \dots, u_N^{(m)} \right) : m = 1, \dots, M \right\} \quad (12)$$

from the joint posterior distribution  $p(\boldsymbol{\theta}, \mathbf{b}_1, \dots, \mathbf{b}_N, u_1, \dots, u_N | \mathbf{y})$  whose margin  $p(\boldsymbol{\theta} | \mathbf{y})$  is the posterior distribution of interest, i.e.,  $p(\boldsymbol{\theta} | \mathbf{y}) \propto L(\boldsymbol{\theta}) p(\boldsymbol{\theta})$ . With respect to the intended clustering, a well-known problem arising from the invariance of the likelihood  $L(\boldsymbol{\theta})$  under permutation of the component labels is resolved using the relabeling algorithm of Stephens (2000) adapted for use in the context of our model.

### 2.3. Clustering procedure

It is explained in Section 1.1 that model-based clustering is based on the estimated values of the individual component probabilities (2). Within the Bayesian framework, the natural estimates of the values  $p_{i,k}$ ,  $i = 1, \dots, N$ ,  $k = 1, \dots, K$ , are their posterior means, MCMC estimates of which are easily obtainable from the generated posterior sample  $\mathcal{S}_M$ , i.e.,

$$\begin{aligned} \hat{p}_{i,k} &= \mathbb{E}\{p_{i,k}(\boldsymbol{\theta}) | \mathbf{Y} = \mathbf{y}\} = \mathbb{P}(U_i = k | \mathbf{Y} = \mathbf{y}) \\ &= \int p_{i,k}(\boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta} \approx \frac{1}{M} \sum_{m=1}^M p_{i,k}(\boldsymbol{\theta}^{(m)}). \end{aligned} \quad (13)$$

Nevertheless, other characteristics of the posterior distributions of the component probabilities, e.g., the posterior medians, can also be used. On top of that, uncertainty in the classification can be to some extent taken into account by exploring either the full posterior distribution of the component probabilities, or by calculating their credible intervals. We illustrate this in Section 3.7.

### 2.4. Selecting a number of mixture components

Selecting a suitable number of mixture components, that is, selecting a number of clusters if this is not known in advance, coincides with the problem of model selection. In the area of mixture models, models with different numbers of components are usually fitted and then compared by a suitable characteristic of model complexity and model fit. Within the **mixAK**

package, two approaches can easily be exploited that we briefly describe in the rest of this section.

### *Penalized expected deviance*

The most commonly used Bayesian characteristic of model complexity and model fit is probably the deviance information criterion (DIC; Spiegelhalter, Best, Carlin, and van der Linde 2002). Nevertheless, as shown by Celeux, Forbes, Robert, and Titterton (2006), its use in mixture models is controversial. An alternative measure, the penalized expected deviance (PED), derived from cross-validation arguments, was suggested by Plummer (2008). It does not suffer from the controversies described by Celeux *et al.* (2006) and can be used even in the context of the mixture models. Also the **mixAK** package exploits the PED for model comparison. The PED, where lower values indicate a better model, is defined as

$$\text{PED} = \mathbf{E}\{D(\boldsymbol{\theta}) \mid \mathbf{Y} = \mathbf{y}\} + p_{opt}. \quad (14)$$

In Equation 14,  $D(\boldsymbol{\theta})$  is the observed data deviance (Equation 11), and  $p_{opt}$  is the penalty term called optimism whose value can be estimated by the use of importance sampling and the use of two parallel chains.

### *Posterior distribution of the deviances*

An alternative procedure for model selection has been suggested by Aitkin, Liu, and Chadwick (2009), later described also in Aitkin (2010, Chapters 7 and 8). They propose basing the model comparison on the full posterior distribution of the deviances. They argue that the model comparison based on one-number criteria (like DIC but also the previously-discussed PED) ignores the uncertainty in the comparison which increases with the increasing number of the model parameters. On the other hand, this uncertainty is taken into account when using the whole posterior distribution of the deviance for the comparison.

Suppose that we want to compare model 1 (with  $K = K_1$ ) and model 2 (with  $K = K_2 > K_1$ ) having the deviances  $D_{K_1}(\boldsymbol{\theta})$  and  $D_{K_2}(\boldsymbol{\theta})$ , respectively. The posterior probability

$$P_{K_2, K_1}(\mathbf{y}) = \mathbf{P}\{D_{K_2}(\boldsymbol{\theta}) - D_{K_1}(\boldsymbol{\theta}) < 0 \mid \mathbf{Y} = \mathbf{y}\} \quad (15)$$

now quantifies our certainty on whether model 2 is better than model 1. The penalization for the increasing complexity of the model is included implicitly in this procedure as models with more parameters lead to a more diffuse posterior distribution of the deviance and hence decrease in  $P_{K_2, K_1}(\mathbf{y})$ . Aitkin *et al.* (2009, Section 4) further suggest calculating the posterior probability

$$P_{K_2, K_1}^*(\mathbf{y}) = \mathbf{P}\{D_{K_2}(\boldsymbol{\theta}) - D_{K_1}(\boldsymbol{\theta}) < -2 \log(9) \doteq -4.39 \mid \mathbf{Y} = \mathbf{y}\}. \quad (16)$$

They argue that if this probability is high (0.9 or more), we have quite strong evidence in favor of a model with  $K = K_2$  over a model with  $K = K_1$ .

Finally, Aitkin (2010) suggests performing an overall comparison of the posterior distributions of deviances by plotting the posterior cumulative distribution functions (CDFs)

$$F_{D_j}(d \mid \mathbf{Y} = \mathbf{y}) = \mathbf{P}\{D_j(\boldsymbol{\theta}) \leq d \mid \mathbf{Y} = \mathbf{y}\}, \quad j = 1, 2. \quad (17)$$

of the deviances. We illustrate this in Section 3.8.



### 3. Example using the package `mixAK`

The aim of this section is to provide a detailed step-by-step R analysis of a particular dataset, to highlight the most important features of package `mixAK` and to show how to extract the most important results. Even though we comment on the results shown in this manuscript in most cases as well, it goes beyond the scope of this paper to provide them in full context or to explain their meaning in detail. All this is given in an accompanying methodological paper (Komárek and Komárková 2013) where the same dataset is analyzed. This section can also be considered as a user manual which allows the readers to run their own similar analyses involving only a mild modification of the example code.

#### 3.1. Data

Longitudinal data clustering capabilities of the R package `mixAK` will be illustrated on the analysis of a subset of the data from a Mayo Clinic trial on 312 patients with primary biliary cirrhosis (PBC) conducted in 1974–1984 (Dickson, Grambsch, Fleming, Fisher, and Langworthy 1989). The primary data are available at <http://lib.stat.cmu.edu/datasets/pbcseq>. In this paper, only  $N = 260$  subjects known to be alive at 910 days of follow-up, and only the longitudinal measurements up to this point will be considered. The corresponding data are available as a `data.frame` `PBC910` in the `mixAK` package. Upon loading the package and the data, we print the columns important for our analysis and few tail rows to describe the data structure.

```
R> library("mixAK")
R> data("PBC910", package = "mixAK")
R> tail(PBC910)[, c(1, 2, 3, 6:9)]
```

	id	day	month	lbili	platelet	spiders	jspiders
913	311	187	6.14	0.405	382	NA	NA
914	311	397	13.04	0.642	408	0	0.266
915	312	0	0.00	1.856	200	1	0.886
916	312	206	6.77	1.705	189	0	0.188
917	312	390	12.81	2.001	148	0	0.173
918	312	775	25.46	2.791	138	1	0.736

It is a longitudinal dataset with one row per visit. There are 1 to 5 visits per subject (identified by column `id`) performed at time of `day` days (`month` months) of follow-up. At each visit, measurements of three markers ( $R = 3$ ) are recorded: continuous logarithmic bilirubin (`lbili`), discrete platelet count (`platelet`) and dichotomous indication of blood vessel malformations (`spiders`). The column `jspiders` is a jittered version of `spiders` which we shall use for drawing some descriptive plots.

As it is exemplified on data of subject `id = 311`, the value of some of the markers considered might be missing at some visits. In this case, the value of the dichotomous `spiders` is not available at the visit performed at 187 days of follow-up. Still, the non-missing values of variables `lbili` and `platelet` from the visit at 187 days shall contribute to the likelihood, the estimation and clustering procedure. We take missingness of any of the outcome variables into account by modifying the expressions for the  $i$ th subject likelihood contributions (Equations 7

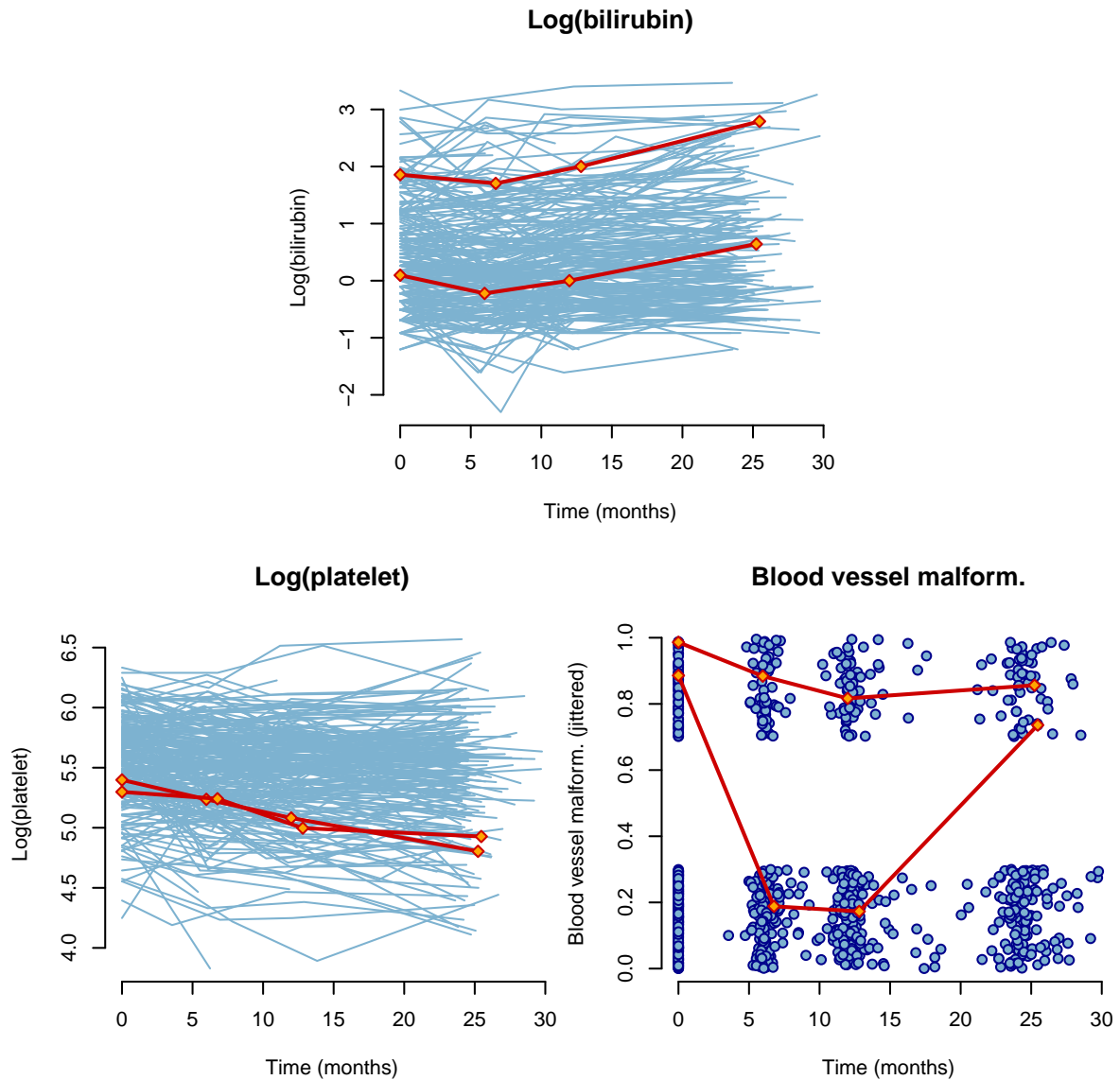


Figure 1: Observed (transformed) longitudinal profiles of considered markers, red lines: profiles of two selected subjects (id 2 and 312).

and 8) such that for given  $r = 1, \dots, R$ , the product over  $j$  takes only the available outcome values into account. Note that this “all available cases” approach is valid as soon as the missingness mechanism can be assumed to be ignorable: data being missing at random or completely at random in the classical taxonomy of Rubin (1976).

In the rest of this section, the random vectors  $\mathbf{Y}_{i,1}, \mathbf{Y}_{i,2}, \mathbf{Y}_{i,3}, i = 1, \dots, N$ , intended for the cluster analysis, shall correspond to the values of `lbili`, `platelet`, `spiders`, respectively. The column `month` shall lead to the time vectors  $\mathbf{t}_1, \dots, \mathbf{t}_N$ . No other covariates will be used in the underlying model.

### 3.2. Basic data exploration

Basic exploration of longitudinal data may consist of drawing longitudinal profiles observed, possibly highlighting the profiles of selected subjects. For this purpose, package **mixAK** offers two functions:

- `getProfiles()` which creates a list of `data.frames` (one `data.frame` per subject) with selected variables;
- `plotProfiles()` which creates a spaghetti graph with observed longitudinal profiles per subject.

As an illustration, we extract the longitudinal profiles of variables `lbili`, `platelet`, `spiders`, `jspiders` and then draw the longitudinal profiles of `lbili` while highlighting the first and the last subject in the dataset (`id` 2 and 312), see the upper panel of Figure 1.

```
R> ip <- getProfiles(t = "month", y = c("lbili", "platelet", "spiders",
+   "jspiders"), id = "id", data = PBC910)
R> plotProfiles(ip = ip, data = PBC910, var = "lbili", tvar = "month",
+   main = "Log(bilirubin)", highlight = c(1, length(ip)),
+   xlab = "Time (months)", ylab = "Log(bilirubin)")
```

The remaining panels of Figure 1 are drawn analogously while using arguments `trans = log` (lower left panel showing the longitudinal profiles of the log-transformed response `platelet`) and `lines = FALSE`, `points = TRUE` (lower right panel showing jittered values of the dichotomous response `spiders`) of the function `plotProfiles()`. The reason for drawing the logarithmic transformation of the variable `platelet` is that a log link will be proposed in Section 3.3 for this response and the created plot then better helps to choose a mean structure of the related GLMM. Further, with a dichotomous longitudinal response, it is almost impossible to draw a fully informative plot of observed data and according to the best of our knowledge, no standard way of doing that in the literature exists. The lower right-hand panel of Figure 1 with vertically jittered values of the response `spiders` is then a possible option which at least allows us into some extent to compare visually the proportions of zeros and ones at each occasion. An alternative could be a suitably adapted spine plot but we do not pursue this option here.

### 3.3. Model

Being partially motivated by Figure 1, by the nature of the considered longitudinal markers and also by other considerations, the following generalized linear mixed models in assumption (A1) of the MMGLMM shall be exploited for the analysis:

- The continuous nature of the variable `lbili` suggests considering a Gaussian GLMM, i.e., a linear mixed model, for the random vectors  $\mathbf{Y}_{i,1}$ ,  $i = 1, \dots, N$ . Further, Figure 1 together with additional exploration of individual longitudinal profiles (not shown) suggest modeling the evolution of each subject by a line over time where the parameters of the line may differ across subjects. This leads to the following random intercept and random slope model for the means of the elements of the response vectors:

$$E(Y_{i,1,j} \mid \mathbf{B}_{i,1} = \mathbf{b}_{i,1}) = b_{i,1,1} + b_{i,1,2}t_{i,j}, \quad \mathbf{b}_{i,1} = (b_{i,1,1}, b_{i,1,2})^\top,$$

$i = 1, \dots, N, j = 1, \dots, n_i$ . Additional refinement of the model, for instance towards capturing possibly non-linear evolution of the outcome over time, is of course possible. Nevertheless, this goes beyond the scope of this software related paper.

In a general notation of Section 2.1, there are no fixed effects, and the random effects covariate vector  $\mathbf{z}_{i,1,j}$  equals  $(1, t_{i,j})^\top$ . The dispersion parameter  $\phi_1$  is an unknown residual variance of the underlying linear mixed model.

- (ii) The count nature of the variable `platelet` leads us to consider a Poisson GLMM for the random vectors  $\mathbf{Y}_{i,2}$ ,  $i = 1, \dots, N$ . Analogously to the previous case of the variable `lbili`, exploration of individual longitudinal profiles (Figure 1 and others) suggests considering subject specific lines over time for logarithmically linked means of the elements of the response vectors:

$$\log\left\{\mathbb{E}(Y_{i,2,j} \mid \mathbf{B}_{i,2} = \mathbf{b}_{i,2})\right\} = b_{i,2,1} + b_{i,2,2}t_{i,j}, \quad \mathbf{b}_{i,2} = (b_{i,2,1}, b_{i,2,2})^\top,$$

$i = 1, \dots, N, j = 1, \dots, n_i$ . Analogously to the model for the variable `lbili`, there are no fixed effects, and the random effects covariate vector  $\mathbf{z}_{i,2,j}$  equals  $(1, t_{i,j})^\top$ . The dispersion parameter  $\phi_2$  is now constant and equal to 1.

- (iii) The dichotomous nature of the variable `spiders` dictates to use a Bernoulli GLMM for which the logit link is a classical choice. The assumed mean structure of the model is the following:

$$\text{logit}\left\{\mathbb{E}(Y_{i,3,j} \mid B_{i,3} = b_{i,3}; \alpha_3)\right\} = b_{i,3} + \alpha_3 t_{i,j},$$

$i = 1, \dots, N, j = 1, \dots, n_i$ . In this case, the fixed effects covariate vector is  $\mathbf{x}_{i,3,j} = (t_{i,j})$ , and the random effects covariate vector is  $\mathbf{z}_{i,3,j} = (1)$ .

Inclusion of the random intercept in the model above is motivated by observing that subjects differ in predisposition towards development of the blood vessel malformations. Due to the fact that we cannot rule out change of this predisposition over time, an additional linear term is included in the model. Nevertheless, due to the relatively low number of repeated measurements per subject, it is difficult to effectively estimate a model with a dichotomous outcome allowing also for subject-specific slopes, which are thus included only as a fixed effect.

In summary, the model parameters (5) common to all clusters are the slope  $\alpha_3$  from the logit model for the variable `spiders`, and the dispersion parameter  $\phi_1$  from the linear mixed model for the variable `lbili`, i.e.,  $\boldsymbol{\xi} = (\alpha_3, \phi_1)^\top$ . The joint random effects vectors  $\mathbf{B}_1, \dots, \mathbf{B}_N$  are five-dimensional,  $\mathbf{B}_i = (B_{i,1,1}, B_{i,1,2}, B_{i,2,1}, B_{i,2,2}, B_{i,3})^\top$ ,  $i = 1, \dots, N$ , with the overall mean  $\boldsymbol{\beta} = (\beta_1, \beta_2, \beta_3, \beta_4, \beta_5)^\top$ , and the overall covariance matrix  $\mathbb{D} = (d_{l,m})_{l,m=1,\dots,5}$ . Initially, a model with  $K = 2$  mixture components will be fitted.

### 3.4. Posterior Markov chain Monte Carlo simulation

Two functions from the `mixAK` package are related to the generation of a posterior sample (12) needed for the inference and the clustering. They include:

- `GLMM_MCMC()` is the main function which runs the MCMC algorithm. Internally, most of the calculation is provided by a compiled C++ code to reduce the computational time. Using the default values of the input arguments, the function generates two MCMC samples started from two sets of initial values which is primarily needed to calculate the penalized expected deviance (PED) useful for model comparison as explained in Section 2.4. The function returns an object which is a `list` and of class `'GLMM_MCMClist'`. It contains some information concerning the model and two objects of class `'GLMM_MCMC'` holding the two sampled chains, initial values used to start the MCMC, specific choices of hyperparameters of the prior distribution and basic posterior summary statistics. Several methods which we shall introduce subsequently are available for objects of classes `'GLMM_MCMClist'` and `'GLMM_MCMC'` to handle and visualize the results.
- `NMixRelabel()` is a generic function with methods for objects of classes `'GLMM_MCMClist'` and `'GLMM_MCMC'` which applies the relabeling algorithm and returns appropriately modified input object. Analogously to the `GLMM_MCMC()` function, the main calculation is internally conducted using the compiled C++ code.

As an illustration, we run the MCMC algorithm for  $(100 \times 10)$  burn-in and  $(1000 \times 10)$  subsequent iterations with 1:10 thinning to get two samples of length  $M = 1000$  from the joint posterior distribution obtained by starting the MCMC from two different sets of initial values. Note that a longer MCMC is usually needed to get reliable results but we keep it rather short here to make the example quickly reproducible. By default, the two chains are generated sequentially. Nevertheless, on multi-core processors this task might be parallelized using the tools provided by the base R package `parallel` by setting the `parallel` argument to `TRUE`.

```
R> set.seed(20042007)
R> mod <- GLMM_MCMC(y = PBC910[, c("lbili", "platelet", "spiders")],
+   dist = c("gaussian", "poisson(log)", "binomial(logit)"),
+   id = PBC910[, "id"], x = list(lbili = "empty", platelet = "empty",
+   spiders = PBC910[, "month"]), z = list(lbili = PBC910[, "month"],
+   platelet = PBC910[, "month"], spiders = "empty"),
+   random.intercept = rep(TRUE, 3), prior.b = list(Kmax = 2),
+   nMCMC = c(burn = 100, keep = 1000, thin = 10, info = 100),
+   parallel = FALSE)
```

The meaning of the most important arguments of the `GLMM_MCMC` function is briefly the following. The argument `y` is a `data.frame` with the observed values of the longitudinal markers in its columns. The mean structure of the GLMM's from assumption (A1) is indicated by arguments `x` (a `list` of vectors/matrices/data frames of the fixed effects covariates except for the intercept term), `z` (a `list` of vectors/matrices/data frames of the random effects covariates except for the intercept term) and a logical argument `random.intercept`. Note that we assume a hierarchically centered GLMM (Gelfand, Sahu, and Carlin 1995) where the random effects have in general non-zero mean (their overall mean  $\beta$  is given by Equation 9). Hence, the marker specific parts of the `lists` in `x` and `z` arguments may not contain the same variables. Otherwise, the model parameters would become unidentifiable as  $\beta$  would be estimated twice, once through the mixture means  $\mu_1, \dots, \mu_K$  and once through the corresponding fixed effects. The keyword `"empty"` is used to indicate that there are no fixed or random effects,

respectively in a model for a specific marker. Purely to simplify the programming work, the `GLMM_MCMC` function is coded such that the intercept is always included in the model and the fact whether it is random or fixed is indicated by argument `random.intercept`. Assumed distributions and the link functions for each marker are indicated by the argument `dist`.

The only obligatory part of the prior distribution which has to be specified by the user is the number of mixture components (clusters) given as a `Kmax` element of the `list` in the `prior.b` argument. All other values of prior hyperparameters are selected automatically to achieve a weakly informative prior distribution using the guidelines described in Komárek and Komárková (2013, Appendix A). The function `GLMM_MCMC` also automatically generates reasonable initial values to start the MCMC simulation. User-defined prior hyperparameters and initial values can be supplied by using the appropriate values of the arguments `prior.alpha`, `prior.b`, `prior.eps`, and `init.alpha`, `init2.alpha`, `init.b`, `init2.b`, `init.eps`, `init2.eps`, respectively.

Before we proceed, we point out that the object `mod` returned by the function `GLMM_MCMC` is a `list` with two main elements `mod[[1]]` and `mod[[2]]` holding the two sampled chains, their initial values, selected posterior summary statistics based on these chains and some additional quantities derived from both sampled chains. The class of the `mod` object is set to `'GLMM_MCMClist'`, the class of the main elements `mod[[1]]` and `mod[[2]]` is `'GLMM_MCMC'`. Basic information concerning the structure of the object `mod` will be given in the subsequent sections. For more details and a description of how to change the default values of the prior hyperparameters and the initial values to start the MCMC simulation, we refer the reader to the package documentation and its vignette.

As it is mentioned in Section 2.2, a problem arising from the invariance of the model likelihood under permutation of the component labels must be resolved and one possibility is to apply a suitable relabeling algorithm. When running the `GLMM_MCMC` function, only a simple relabeling based on the ordering of the first margins of the mixture means was applied. Nevertheless, as illustrated in Stephens (2000, Section 3.1), the relabeling applied above is based in fact on imposing an artificial identifiability constraint and is not always satisfactory. This is especially the case in situations when the chosen identifiability constraint does not separate the mixture components well. Due to this reason, Stephens (2000) suggests a relabeling algorithm which arises from attempting to minimize the posterior expected loss under a class of loss functions suitable for assessment of a clustering procedure. To apply his algorithm on our posterior samples, we run the `NMixRelabel` function (it is done separately for each chain) with its `type` argument set to `"stephens"`:

```
R> mod <- NMixRelabel(mod, type = "stephens", keep.comp.prob = TRUE)
```

Objects `mod[[1]]` and `mod[[2]]` have the same structure as before with the exception that all results which are not invariant towards label switching were (re-)calculated to reflect the new labeling of the mixture components. Finally, by setting `keep.comp.prob = TRUE` in the `NMixRelabel` call, we additionally generated the posterior samples  $p_{i,k}(\boldsymbol{\theta}^{(m)})$ ,  $i = 1, \dots, N$ ,  $k = 1, \dots, K$ ,  $m = 1, \dots, M$  (Equation 13), that subject  $i$  belongs to the  $k$ th group (where  $k$  refers to a new labeling of the components) which will be used in Section 3.7 for the purposes of clustering.

### 3.5. Posterior samples and basic convergence diagnostics

Posterior samples of the primary model parameters, random hyperparameters and some additional derived quantities are kept in the `mod[[1]]` and `mod[[2]]` objects as matrix or vector elements. The most important ones include: `Deviance` (sampled model deviance  $L(\boldsymbol{\theta})$ , Equation 11), `alpha` (sampled fixed effects  $\alpha_1, \dots, \alpha_R$ ), `sigma_eps` (sampled square roots of dispersion parameters  $\phi_1, \dots, \phi_R$ ), `mixture_b` (sampled overall random effects means  $\boldsymbol{\beta}$ , Equation 9, and standard deviations and correlations derived from the overall random effects covariance matrix  $\mathbb{D}$ , Equation 10). As illustration, we print the first ten sampled values of the model deviance  $L(\boldsymbol{\theta})$  from the first chain:

```
R> print(mod[[1]]$Deviance[1:10], digits = 9)

[1] 14086.4438 14095.8671 14113.7103 14097.6683 14109.1802 14102.2894
[7] 14110.1587 14098.9003 14109.3774 14096.8390
```

Classical tools for convergence diagnostics as implemented, e.g., in the R package `coda` (Plummer, Best, Cowles, and Vines 2006) can be used to evaluate the convergence of the performed MCMC simulation. As an illustration, we use the `coda` package routine `autocorr()` to calculate estimated autocorrelations (in the first and the second chain) in our Markov chain of the model deviances (the output was edited into two columns):

```
R> library("coda")
R> DevChains <- mcmc.list(mcmc(mod[[1]]$Deviance), mcmc(mod[[2]]$Deviance))
R> autocorr(DevChains)
```

[[1]]	[[2]]
, , 1	, , 1
Lag 0 1.00000	Lag 0 1.00000
Lag 1 0.24155	Lag 1 0.21846
Lag 5 0.10663	Lag 5 0.12404
Lag 10 0.06345	Lag 10 0.03833
Lag 50 0.00043	Lag 50 -0.01368

Two extra routines are available within the `mixAK` package to access the posterior samples and produce basic diagnostics plots. These are:

- `NMixChainComp()` is a generic function with a method for objects of class ‘GLMM\_MCMC’ which extracts the posterior samples of the mixture parameters from the resulting object: weights  $\mathbf{w}$ , means  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$ , covariance matrices  $\mathbb{D}_1, \dots, \mathbb{D}_K$  or their derivatives (standard deviations and correlations based on  $\mathbb{D}_1, \dots, \mathbb{D}_K$ ). It has a logical argument `relabel` that determines whether the original sample is required (`relabel = FALSE`) or a relabeled sample (with default `relabel = TRUE`) as calculated by the `NMixRelabel` function in Section 3.4.
- `tracePlots()` is again a generic function with methods for objects of classes ‘GLMM\_MCMC’ and ‘GLMM\_MCMClist’ which produces trace plots (parallel ones if applied to the object

of class ‘GLMM\_MCMClist’) of a selected class of model parameters. In case that trace plots of mixture related parameters are required, a logical argument `relabel` with the same meaning as in the case of the `NMixChainComp()` function determines whether the original or the relabeled sample is to be plotted.

We first illustrate the use of the `NMixChainComp()` function and extract from the object `mod[[1]]` (the first sampled chain) the relabeled sample of the mixture means  $\mu_1, \dots, \mu_K$  and then print the first three sampled values.

```
R> muSamp1 <- NMixChainComp(mod[[1]], relabel = TRUE, param = "mu_b")
R> print(muSamp1[1:3, ])

      mu1.1  mu1.2 mu1.3  mu1.4 mu1.5 mu2.1  mu2.2 mu2.3  mu2.4
[1,] -0.258 0.00645  5.59 -0.00690 -3.82 1.240 0.00322  5.43 -0.00729
[2,] -0.257 0.00345  5.59 -0.00534 -4.41 0.951 0.01426  5.39 -0.00277
[3,] -0.323 0.00625  5.61 -0.00583 -3.83 0.838 0.01144  5.41 -0.00315
      mu2.5
[1,] -1.101
[2,] -1.069
[3,] -0.563
```

The first five columns of matrix `muSamp1` contain the sampled values of  $\mu_1$ , the remaining five columns contain the sampled values of  $\mu_2$ . By changing the value of the argument `param`, sampled values of other mixture parameters are provided. In particular, `param` values of `"w_b"`, `"var_b"`, `"sd_b"`, `"cor_b"`, provide the sampled values of mixture weights, variances (diagonal elements of matrices  $\mathbb{D}_1, \dots, \mathbb{D}_K$ ), their square roots (standard deviations), and correlation coefficients derived from matrices  $\mathbb{D}_1, \dots, \mathbb{D}_K$ , respectively. Further, the `param` values of `"Sigma_b"`, `"Q_b"`, `"Li_b"` provide the sampled values of the lower triangles of the mixture covariance matrices  $\mathbb{D}_1, \dots, \mathbb{D}_K$ , their inversions, and Cholesky factors of their inversions, respectively.

Trace plots are a useful basic tool for exploring the behavior of the sampled Markov chains. When two parallel chains are generated, it is useful to draw both chains in two different colors into one plot as it is done by the `mixAK` function `tracePlots`. As illustration, we draw the trace plots of the two parallel chains of the model deviance  $D(\theta)$ , see Figure 2.

```
R> tracePlots(mod, param = "Deviance")
```

By changing the argument `param`, trace plots of other model parameters are drawn. In particular, by setting the argument `param` to `"alpha"` and `"sigma_eps"`, respectively, trace plots of the fixed effects vector  $\alpha_1, \dots, \alpha_R$  and the square roots of the dispersion parameters  $\phi_1, \dots, \phi_R$ , respectively, are drawn. Likewise, the `param` argument values of `"Eb"`, `"SDb"` and `"Corb"`, respectively, lead to trace plots for the overall means  $\beta$  (Equation 9) of the random effects, and the standard deviations and correlation coefficients, respectively, derived from the overall covariance matrix  $\mathbb{D}$  (Equation 10).

Further, the `param` argument values of `"w_b"`, `"mu_b"` and `"sd_b"` can be used to draw the trace plots of the mixture weights  $w_1, \dots, w_K$ , mixture means  $\mu_1, \dots, \mu_K$  and standard deviations derived from the mixture covariance matrices  $\mathbb{D}_1, \dots, \mathbb{D}_K$  where the user can also



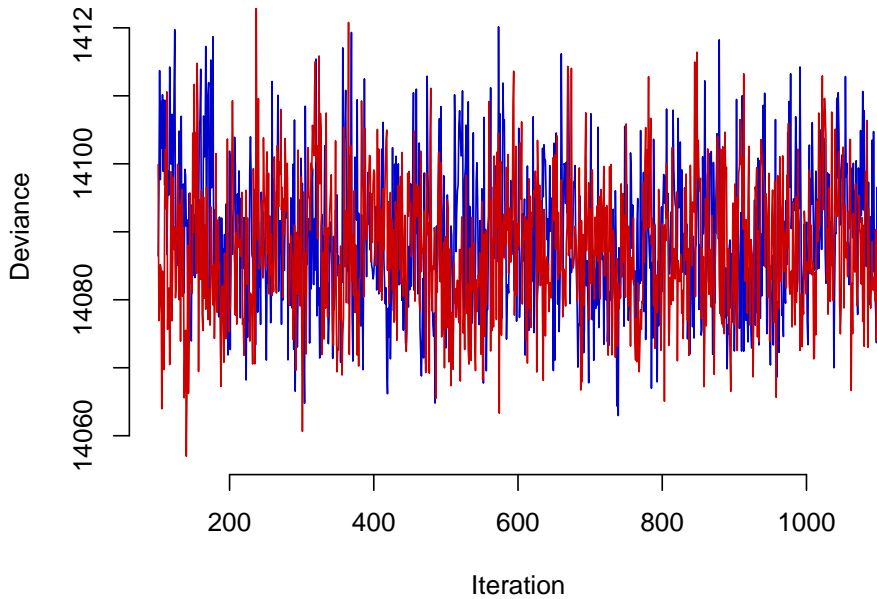


Figure 2: Trace plots of the two parallel chains of the model deviance  $D(\boldsymbol{\theta})$ .

decide (by setting the argument `relabel`) whether the trace plots of the mixture parameters are to be drawn as sampled or after relabeling as performed in Section 3.4. For example, the trace plots of relabeled mixture means  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$  are shown in Figure 3 generated by the following code:

```
R> tracePlots(mod, param = "mu_b", relabel = TRUE)
```

Finally, when the first argument of the `tracePlots` function is changed to `mod[[1]]` or `mod[[2]]` only the single trace plots of the first and the second chain, respectively, are drawn.

### 3.6. Posterior summary statistics

To summarize the estimated models, we usually calculate the posterior summary statistics and credible intervals for important model parameters. In the context of the MMGLMM, these include: the fixed effects  $\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_R$  and the dispersion parameters  $\phi_1, \dots, \phi_R$ , and the overall mean  $\boldsymbol{\beta}$  (Equation 9) and the overall covariance matrix  $\mathbb{D}$  (Equation 10) of the random effects. Furthermore, the posterior summary statistics of the observed data deviance  $D(\boldsymbol{\theta})$  (Equation 11) is usually reported to provide a basic model fit characteristic. Note that up to now mentioned quantities are invariant towards label switching and hence their posterior summary statistics might be calculated even without performing any relabeling. Having the subsequent clustering in mind, posterior summary statistics for the mixture parameters (weights  $w_1, \dots, w_K$ , means  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$ , covariance matrices  $\mathbb{D}_1, \dots, \mathbb{D}_K$ ) are of additional interest. Nevertheless, a suitable relabeling of the posterior sample must be first conducted to calculate these, as we did in Section 3.4.

In this section, we first show how to easily obtain basic posterior summary statistics of the important model parameters by the mean of routines implemented in the **mixAK** package. Second, we exemplify usage of the posterior samples stored in the object `mod` created in Section 3.4 in connection with the **coda** package to obtain more detailed posterior summaries.

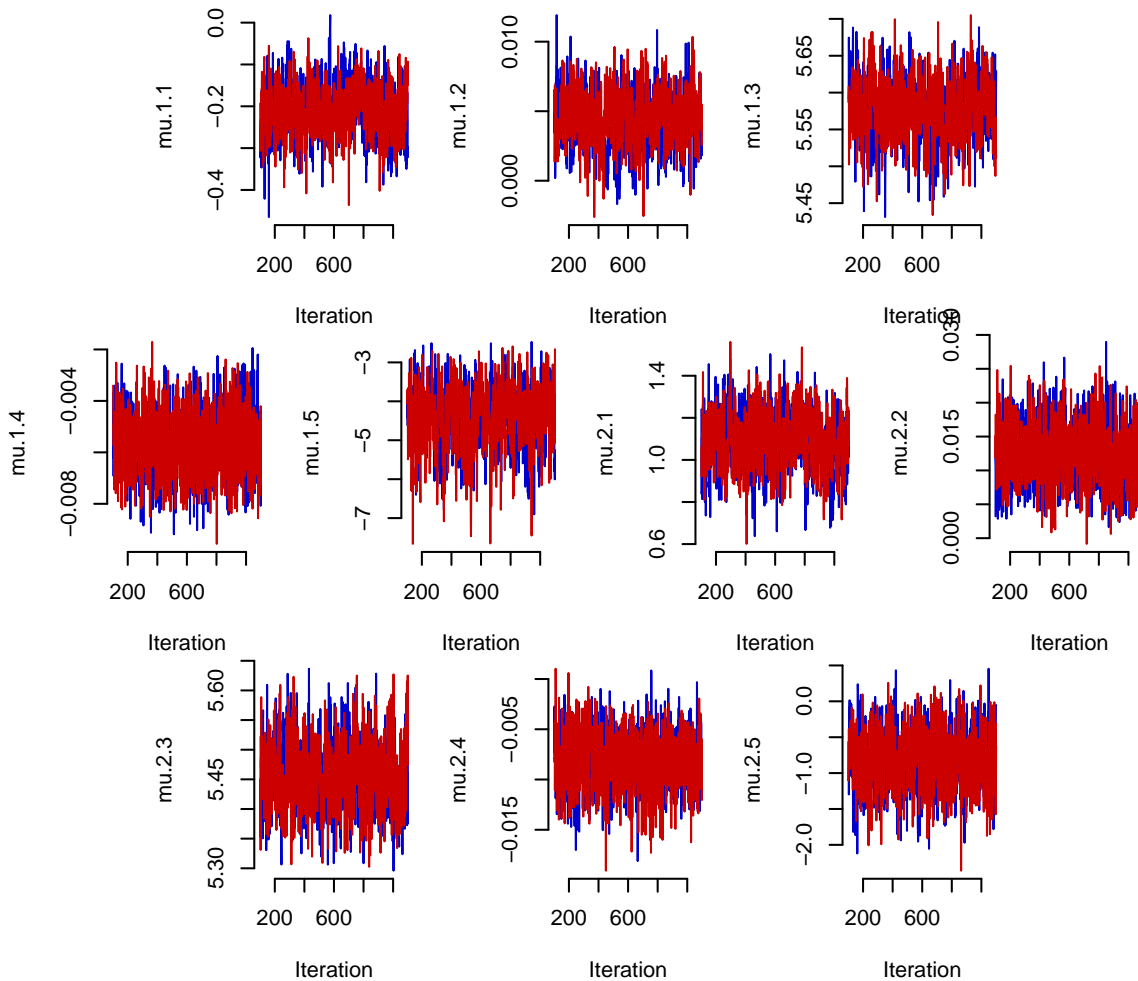


Figure 3: Trace plots of the elements of the mixture means  $\mu_1, \dots, \mu_K$  after Stephens' relabeling.

#### *Posterior summary statistics for the GLMM related parameters*

Basic posterior summary statistics for many of the above-mentioned quantities have in fact already been calculated by the function `GLMM_MCMC` and are stored as the following components of the objects `mod[[1]]` and `mod[[2]]`: `summ.Deviance` (observed data deviance  $D(\theta)$ ), `summ.alpha` (fixed effect  $\alpha_1, \dots, \alpha_R$ ), `summ.sigma_eps` (square roots of the dispersion parameters  $\phi_1, \dots, \phi_R$ ), `summ.b.Mean` (overall means  $\beta$  of random effects), `summ.b.SDCorr` (standard deviations and the correlation coefficients derived from the overall covariance matrix  $\mathbb{D}$  of random effects). To inspect their values in a synoptic form, we simply print the object `mod` (output has been shortened).

```
R> print(mod)
```

```
Generalized linear mixed model for 3 responses estimated using MCMC
```

```
=====
```

Penalized expected deviance:

```
-----
D.expect  p(opt)      PED  wp(opt)      wPED
14088.3   74.5  14162.8   74.8  14163.1
```

Deviance posterior summary statistics:

```
-----
          Mean Std.Dev.  Min.  2.5% 1st Qu. Median 3rd Qu. 97.5%  Max.
Chain 1 14089      10.5 14063 14071  14081  14088  14096 14111 14120
Chain 2 14088      10.1 14057 14069  14080  14087  14094 14108 14125
```

Posterior summary statistics for fixed effects:

```
-----
          Mean Std.Dev.  Min.  2.5% 1st Qu. Median 3rd Qu. 97.5%  Max.
Chain 1 0.0277   0.0128 -0.0124 0.0029  0.0189 0.0278  0.0367 0.0519 0.0735
Chain 2 0.0281   0.0130 -0.0209 0.0024  0.0200 0.0282  0.0367 0.0519 0.0769
```

Distribution of random effects is a normal mixture with 2 components

-----  
Posterior summary statistics for moments of mixture for random effects:  
-----

Means:

```

          b.Mean.1(Chain 1) b.Mean.1(Chain 2) b.Mean.2(Chain 1)
Mean          0.3185          0.3103          0.00775
Std.Dev.      0.0557          0.0566          0.00182
Min.          0.1422          0.1130          0.00180
2.5%          0.2056          0.2017          0.00443
1st Qu.       0.2805          0.2716          0.00650
Median        0.3194          0.3097          0.00776
3rd Qu.       0.3545          0.3498          0.00898
97.5%        0.4288          0.4171          0.01151
Max.          0.4996          0.4833          0.01401
          b.Mean.2(Chain 2) b.Mean.3(Chain 1) b.Mean.3(Chain 2)
Mean          0.00785          5.5262          5.5273
...

```

Standard deviations and correlations:

```

          b.SD.1(Chain 1) b.SD.1(Chain 2) b.Corr.2.1(Chain 1)
Mean          0.8736          0.8728          0.0527
...

```

Posterior summary statistics for standard deviations  
of residuals of continuous responses:

```
-----
          Mean Std.Dev.  Min.  2.5% 1st Qu. Median 3rd Qu. 97.5%  Max.
Chain 1 0.314   0.0101 0.282 0.294  0.307  0.314  0.321 0.334 0.344
Chain 2 0.314   0.0104 0.284 0.293  0.307  0.314  0.321 0.337 0.352
```

On top of the output, we have a summary of the PED introduced in Section 2.4. We return to it in Section 3.8 when discussing the selection of the suitable number of components for this particular example. The rest of the output is devoted to the posterior summary statistics (posterior mean, posterior standard deviation, posterior 0%, 2.5%, 25%, 50%, 75%, 97.5% and 100% quantiles) of the most important model parameters. Among other things, we also directly see the 95% equal-tail credible intervals. In particular, the output shows the posterior summaries of the following model parameters: (i) the model deviance  $D(\boldsymbol{\theta})$ , e.g., its posterior means estimated using the first and the second chain are 14 089 and 14 088, respectively, the estimated posterior medians have values of 14 088 and 14 087, respectively; (ii) the fixed effects  $\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_R$ ; (iii) the overall mean vector  $\boldsymbol{\beta}$  of the random effects; (iv) standard deviations and correlations derived from the overall covariance matrix  $\mathbb{D}$  of the random effects; (v) square roots of the dispersion parameters  $\phi_1, \dots, \phi_R$ .

To get the posterior summary including estimates of Monte Carlo errors in the estimation of the posterior means, one may use the `summary` procedure from the `coda` package applied directly to matrices of sampled values. As illustration, we create a matrix containing the sampled values of the means  $\boldsymbol{\beta}$  of random effects, the fixed effect  $\alpha_3$  (slope from the logit model for a variable `spiders`), and the square root of the dispersion parameter  $\phi_1$  (residual standard deviation from the linear mixed model for a variable `lbili`), and calculate the posterior summary using the `coda` package (output was shortened).

```
R> name.Eb <- paste("b.Mean.", 1:5, sep = "")
R> RegrChain1 <- with(mod[[1]],
+   cbind(mixture_b[, name.Eb], alpha, sigma_eps))
R> RegrChain2 <- with(mod[[2]],
+   cbind(mixture_b[, name.Eb], alpha, sigma_eps))
R> colnames(RegrChain1) <- colnames(RegrChain2) <-
+   c(paste(rep(c("lbili", "platelet", "spiders"), each = 2), ":"),
+     rep(c("Intcpt", "Slope"), 3), sep = ""), "lbili:res_std_dev")
R> summary(mcmc.list(mcmc(RegrChain1), mcmc(RegrChain2)))
```

```
Iterations = 101:1100
Thinning interval = 1
Number of chains = 2
Sample size per chain = 1000
```

1. Empirical mean and standard deviation for each variable, plus standard error of the mean:

	Mean	SD	Naive SE	Time-series SE
lbili:Intcpt	0.31439	0.05628	1.26e-03	1.23e-03
lbili:Slope	0.00780	0.00180	4.03e-05	4.49e-05
platelet:Intcpt	5.52674	0.02215	4.95e-04	4.29e-04
...				

2. Quantiles for each variable:

	2.5%	25%	50%	75%	97.5%
--	------	-----	-----	-----	-------

```

lbili:Intcpt      0.20293  0.27519  0.31443  0.35174  0.42221
lbili:Slope      0.00429  0.00659  0.00782  0.00895  0.01140
platelet:Intcpt  5.48533  5.51110  5.52607  5.54155  5.56978
...

```

The above output is now based on values from both sampled Markov chains. Analogously, other **coda** routines like `HPDinterval()`, `densplot()` and others might be used to calculate the highest posterior density (HPD) credible intervals, the posterior densities and other posterior quantities.

### *Posterior summary statistics for the mixture components*

With respect to the clustering intended, posterior summary statistics of the mixture parameters (weights  $\mathbf{w}$ , means  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$  and covariance matrices  $\mathbb{D}_1, \dots, \mathbb{D}_K$ ) are of primary interest as they provide characteristics of the mixture and hence also of the clusters. To this end, the **mixAK** package offers a function `NMixSummComp()` which reports posterior means of the mixture weights, mixture means and mixture covariance matrices based on the relabeled sample. On top of that, standard deviations and correlations derived from the posterior means of the covariance matrices are reported (output shortened):

```
R> NMixSummComp(mod[[1]])
```

Component 1

```

Weight:  0.589
Mean:    -0.211  0.00426  5.58 -0.00559 -4.27

```

Covariance matrix:

	m1	m2	m3	m4	m5
m1	1.80e-01	8.75e-05	-3.62e-02	-3.71e-04	0.528122
m2	8.75e-05	6.60e-05	9.01e-05	-1.84e-05	0.003039
m3	-3.62e-02	9.01e-05	9.41e-02	-1.31e-04	-0.049586
m4	-3.71e-04	-1.84e-05	-1.31e-04	1.09e-04	0.000837
m5	5.28e-01	3.04e-03	-4.96e-02	8.37e-04	14.986695

Standard deviations: 0.424 0.00813 0.307 0.0105 3.87

Correlation matrix:

	m1	m2	m3	m4	m5
m1	1.0000	0.0254	-0.2780	-0.0835	0.3215
m2	0.0254	1.0000	0.0361	-0.2161	0.0966
m3	-0.2780	0.0361	1.0000	-0.0409	-0.0418
m4	-0.0835	-0.2161	-0.0409	1.0000	0.0207
m5	0.3215	0.0966	-0.0418	0.0207	1.0000

-----

Component 2

```
Weight:  0.411
```

...

That is, for example the posterior means of the parameters of the first mixture component based on the performed MCMC simulation and a relabeled sample (the first chain) are

$$\begin{aligned}\hat{w}_1 &= \mathbb{E}(w_1 \mid \mathbf{Y} = \mathbf{y}) = 0.589, \\ \hat{\boldsymbol{\mu}}_1 &= \mathbb{E}(\boldsymbol{\mu}_1 \mid \mathbf{Y} = \mathbf{y}) = (-0.211, 0.00426, 5.58, -0.00559, -4.27)^\top,\end{aligned}\tag{18}$$

$$\hat{\mathbb{D}}_1 = \mathbb{E}(\mathbb{D}_1 \mid \mathbf{Y} = \mathbf{y})$$

$$= \begin{pmatrix} 0.180 & 0.0000875 & -0.0362 & -0.000371 & 0.528 \\ 0.0000875 & 0.0000660 & 0.0000901 & -0.0000184 & 0.00304 \\ -0.0362 & 0.0000901 & 0.0941 & -0.000131 & -0.0496 \\ -0.000371 & -0.0000184 & -0.000131 & 0.000109 & 0.000837 \\ 0.528 & 0.00304 & -0.0496 & 0.000837 & 15.0 \end{pmatrix}.$$

Analogously to the case of the GLMM related parameters, more detailed posterior summary of the mixture parameters can be obtained by using the `coda` routines with the posterior samples of the mixture parameters extracted by the mean of the `mixAK` function `NMixChainComp()`. For example, more detailed posterior summary statistics and the 95% HPD credible intervals for the mixture means based on the first sampled chain are obtained using (output was shortened):

```
R> summary(mcmc(muSamp1))
```

```
Iterations = 1:1000
Thinning interval = 1
Number of chains = 1
Sample size per chain = 1000
```

1. Empirical mean and standard deviation for each variable, plus standard error of the mean:

	Mean	SD	Naive SE	Time-series SE
mu1.1	-0.21126	0.06511	2.06e-03	5.02e-03
mu1.2	0.00426	0.00197	6.21e-05	1.39e-04
mu1.3	5.57800	0.04174	1.32e-03	4.11e-03
...				

2. Quantiles for each variable:

	2.5%	25%	50%	75%	97.5%
mu1.1	-0.334236	-0.25434	-0.21351	-0.16864	-0.08649
mu1.2	0.000298	0.00292	0.00422	0.00568	0.00790
mu1.3	5.494958	5.55050	5.58047	5.60731	5.65415
...					

```
R> HPDinterval(mcmc(muSamp1))
```

```

      lower  upper
mu1.1 -0.340065 -0.09240
mu1.2  0.000785  0.00825
mu1.3  5.497083  5.65449
...
attr(,"Probability")
[1] 0.95

```

### 3.7. Clustering

In this section, we first discuss the possibilities of characterization of the mixture components in more detail, i.e., clusters that were found by our analysis. Second, we exemplify usage of the posterior samples stored in the object `mod` for classification of individual subjects into those clusters.

#### *Cluster specific mean longitudinal profiles*

As we mentioned in Section 3.5, clusters in our model are characterized by the mixture parameters (weights  $w_1, \dots, w_K$ , means  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$ , covariance matrices  $\mathbb{D}_1, \dots, \mathbb{D}_K$ ) and their derivatives for which we already calculated posterior summary statistics, see, e.g., Equation 18 for the posterior means of the parameters of the first mixture component. To visualize these results and to see better the characteristics of the different clusters, we can calculate estimated values of the cluster specific mean longitudinal profiles of the response variables. To this end, the generic function `fitted()` was extended in the `mixAK` package by a method for objects of class ‘GLMM\_MCMC’. It provides the empirical Bayes estimates of those longitudinal profiles. In the following, let  $\hat{\boldsymbol{\mu}}_k, \hat{\mathbb{D}}_k, k = 1, \dots, K$ , be the posterior means of mixture means and covariance matrices, respectively. Further, let  $\hat{\boldsymbol{\mu}}_{k,r}$  and  $\hat{\mathbb{D}}_{k,r}$  denote for each  $k = 1, \dots, K, r = 1, \dots, R$  a proper subvector of  $\hat{\boldsymbol{\mu}}_k$  and a submatrix of  $\hat{\mathbb{D}}_k$ , respectively, corresponding to the random effects vector for the  $r$ th response variable. In particular, the function `fitted` calculates the values of

$$\begin{aligned}
 & \mathbb{E}(Y_{new,r,j} \mid U_{new} = k; \hat{\boldsymbol{\theta}}) \\
 &= \int \mathbb{E}(Y_{new,r,j} \mid \mathbf{B}_{new,r} = \mathbf{b}_{new,r}; \hat{\boldsymbol{\alpha}}) \varphi(\mathbf{b}_{new,r}; \hat{\boldsymbol{\mu}}_{k,r}, \hat{\mathbb{D}}_{k,r}) d\mathbf{b}_{new,r}, \quad (19)
 \end{aligned}$$

$r = 1, \dots, R, k = 1, \dots, K, j = 1, \dots, n_{new}$  for selected combinations of covariates entering the model expression (3) of  $\mathbb{E}(Y_{new,r,j} \mid \mathbf{B}_{new,r} = \mathbf{b}_{new,r}; \hat{\boldsymbol{\alpha}})$ .

In our example, the only covariates included in the model are the visit times. The following code then calculates the values of (19) for  $r = 1, \dots, R, k = 1, \dots, K$ , and the covariate values  $t_{new,1} = 0, t_{new,2} = 0.3, \dots, t_{new,101} = 30$  (i.e.,  $n_{new} = 101$ ):

```

R> delta <- 0.3
R> tpred <- seq(0, 30, by = delta)
R> fit <- fitted(mod[[1]], x = list("empty", "empty", tpred),
+   z = list(tpred, tpred, "empty"), glmer = TRUE)
R> names(fit) <- c("lbili", "platelet", "spiders")

```

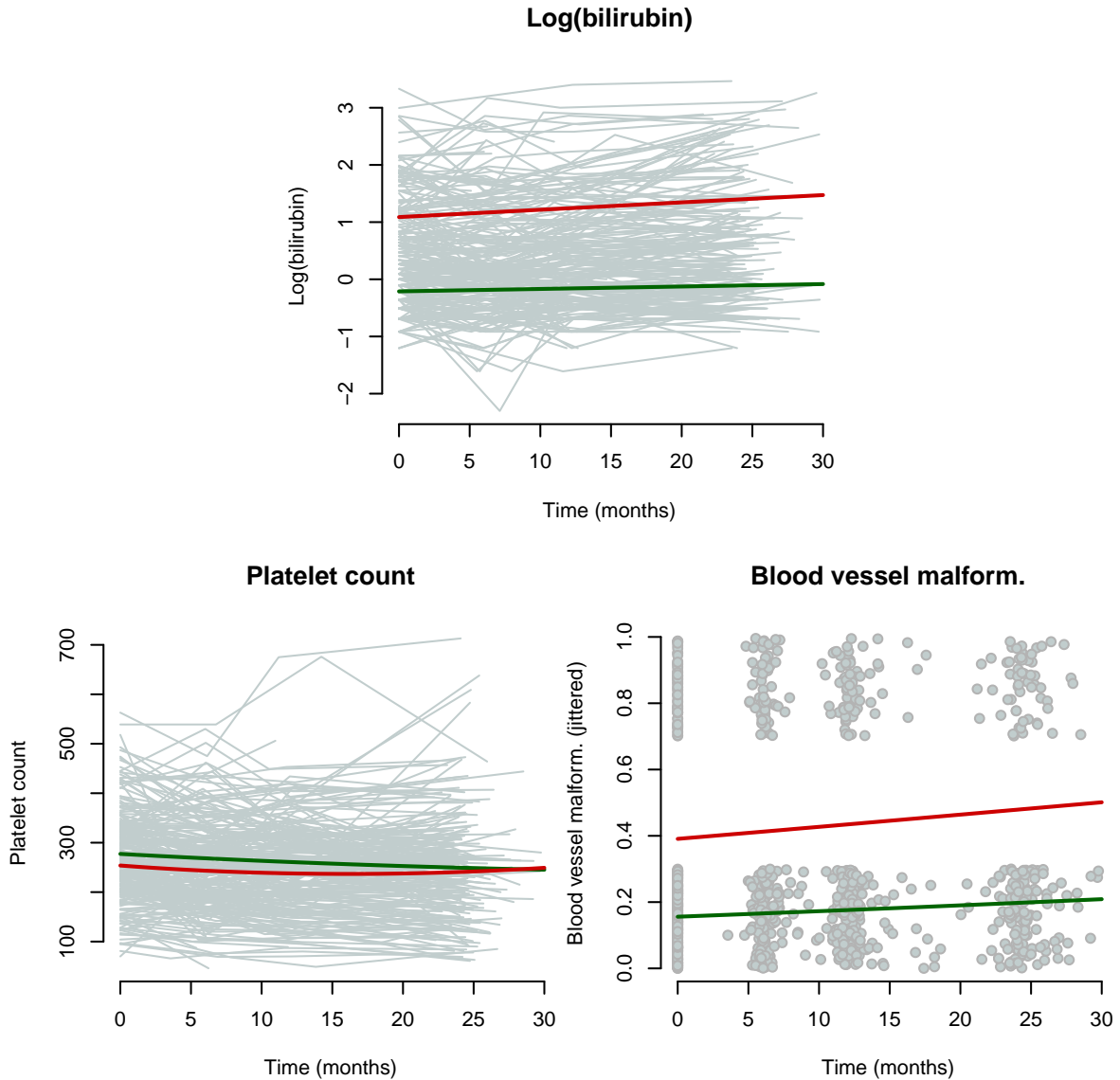


Figure 4: Observed longitudinal profiles of considered markers together with the estimated cluster specific mean profiles (cluster 1 in green, cluster 2 in red).

The covariate combinations for which (19) is to be calculated are given as arguments  $\mathbf{x}$  (fixed effects covariates) and  $\mathbf{z}$  (random effects covariates).

Further, we point out that with `glmer = FALSE` in the call to `fitted`, calculation is faster but the integral in (19) is only rather inaccurately approximated by  $E(Y_{new,r,j} | \mathbf{B}_{new,r} = \hat{\boldsymbol{\mu}}_{k,r}; \hat{\boldsymbol{\alpha}})$ , that is, the random effect values are replaced by their means rather than being integrated out.

The resulting object `fit` is a list of length  $R = 3$ , where each list component is a  $n_{new} \times K$  matrix. The  $k$ th column,  $k = 1, \dots, K$ , of the  $r$ th,  $r = 1, \dots, R$ , matrix gives the calculated values of  $E(Y_{.,r.} | U. = k, \hat{\boldsymbol{\theta}})$  for the covariate combinations corresponding to the vectors or matrices specified by the `x`, `x2`, `z`, `z2` arguments. For instance, the first three values of the



component specific estimated longitudinal profiles for the variable `platelet` are as follows.

```
R> print(fit[["platelet"]][1:3, ], digits = 5)

      [,1] [,2]
[1,] 277.28 253.83
[2,] 276.81 253.18
[3,] 276.34 252.55
```

With respect to the interpretation of the clusters found, it is probably more useful to draw the estimated cluster specific mean longitudinal profiles, perhaps together with the data observed. To achieve this, we use a mild modification of the code shown in Section 3.2 for Figure 1. A plot of observed longitudinal profiles together with the cluster specific mean profiles for a variable `lbili` (see the upper panel of Figure 4) is drawn using the following commands:

```
R> K <- mod[[1]]$prior.b$Kmax
R> c1COL <- c("darkgreen", "red3")
R> plotProfiles(ip = ip, data = PBC910, var = "lbili", tvar = "month",
+   col = "azure3", main = "Log(bilirubin)", xlab = "Time (months)",
+   ylab = "Log(bilirubin)")
R> for (k in 1:K) lines(tpred, fit[["lbili"]][, k], col = c1COL[k], lwd = 2)
```

The remaining panels of Figure 4 are drawn analogously. We now see from Figure 4 that the first (green) cluster is characterized by lower bilirubin values and also by less frequent blood vessel malformations whose probability only slightly rises over time. On the other hand, the second (red) cluster corresponds to subjects with higher bilirubin values and more frequent blood vessel malformations whose occurrence increases more steeply over time. With respect to the longitudinal evolution of the platelet counts, both groups behave almost equally on average.

### *Individual component probabilities*

We explained in Section 1.1 that model-based clustering is based on the individual component probabilities (ICPs)  $p_{i,k} = p_{i,k}(\boldsymbol{\theta}) = \mathbb{P}(U_i = k \mid \mathbf{Y}_i = \mathbf{y}_i; \boldsymbol{\theta})$ ,  $i = 1, \dots, N$ ,  $k = 1, \dots, K$ , see Equation 2. Upon running the `NMIXRELabel()` procedure in Section 3.4, an MCMC sample from the posterior distribution of  $p_{i,k}$  is available within the resulting ‘GLMM\_MCMC’ objects `mod[[1]]` and `mod[[2]]` for each  $i = 1, \dots, N$ ,  $k = 1, \dots, K$ . These are the values

$$p_{i,k}^{(m)} = p_{i,k}(\boldsymbol{\theta}^{(m)}) = \mathbb{P}(U_i = k \mid \mathbf{Y}_i = \mathbf{y}_i; \boldsymbol{\theta}^{(m)}), \quad m = 1, \dots, M, \quad (20)$$

where  $\boldsymbol{\theta}^{(1)}, \dots, \boldsymbol{\theta}^{(M)}$  is an MCMC sample from the posterior distribution  $p(\boldsymbol{\theta} \mid \mathbf{y})$ . In the following paragraphs, we show how to use different characteristics of the posterior distributions of each  $p_{i,k}$  for clustering purposes.

Before we proceed to classification itself, we show how to access the values of (20). They are stored as an  $M \times (N \cdot K)$  matrix `comp.prob` of ‘GLMM\_MCMC’ objects where each row corresponds to one MCMC iteration and the columns of the  $m$ th row provide in a sequence the values of  $p_{1,1}^{(m)}, \dots, p_{1,K}^{(m)}, \dots, p_{N,1}^{(m)}, \dots, p_{N,K}^{(m)}$ . The labeling of components corresponds to that obtained by applying a relabeling algorithm in Section 3.4. For example, the first three sampled values (from the first chain) of the ICPs of the first four subjects in the dataset are as follows:

```
R> print(mod[[1]]$comp.prob[1:3, 1:8])

      P(1,1) P(1,2) P(2,1) P(2,2) P(3,1) P(3,2) P(4,1) P(4,2)
[1,]  0.803  0.197  0.704  0.296  0.84183  0.158  0.054087  0.946
[2,]  0.729  0.271  0.437  0.563  0.16346  0.837  0.013635  0.986
[3,]  0.187  0.813  0.173  0.827  0.00965  0.990  0.000235  1.000
```

Finally, we point out that a similar matrix called `comp.prob_b` is present inside the ‘GLMM\_MCMC’ objects. It provides the values of  $P(U_i = k \mid \mathbf{Y}_i = \mathbf{y}_i; \mathbf{B}_i = \mathbf{b}_i^{(m)}, \boldsymbol{\theta}^{(m)})$  (which are different from those given by Equation 20),  $i = 1, \dots, N$ ,  $k = 1, \dots, K$ ,  $m = 1, \dots, M$ , where  $\mathbf{b}_i^{(m)}$  are the sampled values of the  $i$ th subject random effects. The values of  $P(U_i = k \mid \mathbf{Y}_i = \mathbf{y}_i; \mathbf{B}_i = \mathbf{b}_i^{(m)}, \boldsymbol{\theta}^{(m)})$  might be of interest in some situations. Nevertheless, we shall not use them for any purposes in this manuscript.

### Clustering procedure

As indicated in Section 2.3, the clustering procedure in a Bayesian setting can straightforwardly be based on the posterior means  $\hat{p}_{i,k} = E\{p_{i,k}(\boldsymbol{\theta}) \mid \mathbf{Y} = \mathbf{y}\}$ ,  $i = 1, \dots, N$ ,  $k = 1, \dots, K$ , of the ICPs, see also Equation 13. Their MCMC based estimates are included in the ‘GLMM\_MCMC’ objects as an  $N \times K$  matrix `poster.comp.prob`. We print its first three rows which provide the MCMC estimated values of  $\hat{p}_{i,k}$  for  $i = 1, 2, 3$ ,  $k = 1, 2$ :

```
R> print(mod[[1]]$poster.comp.prob[1:3, ])

      [,1] [,2]
[1,] 0.683 0.317
[2,] 0.558 0.442
[3,] 0.489 0.511
```

The most common classification procedure which assigns subject  $i$  into group  $g(i)$  satisfying  $\hat{p}_{i,g(i)} = \max_{k=1,\dots,K} \hat{p}_{i,k}$  is then achieved by the following code which also stores the values of  $\hat{p}_{i,g(i)}$ ,  $i = 1, \dots, N$ , in a vector `pMean`:

```
R> groupMean <- apply(mod[[1]]$poster.comp.prob, 1, which.max)
R> pMean <- apply(mod[[1]]$poster.comp.prob, 1, max)
R> table(groupMean)
```

```
groupMean
 1    2
161  99
```

That is, 161 and 99 subjects are classified in the first and the second group, respectively, and represented in Figure 4 by the green and the red cluster specific mean profiles, respectively.

Nevertheless, the posterior mean is only one possible one-number characteristic of the posterior distribution. A reasonable, and often even more suitable characteristic is the posterior median which can also be used for the clustering procedure. The posterior medians, together with the 2.5% and 97.5% quantiles of the posterior distributions of the ICPs are available in the `quant.comp.prob` elements of the ‘GLMM\_MCMC’ objects. As illustration, we again provide the posterior medians of the ICPs for the first three subjects:

```
R> print(mod[[1]]$quant.comp.prob[["50%"]][1:3, ])
```

```
      [,1] [,2]
[1,] 0.786 0.214
[2,] 0.615 0.385
[3,] 0.495 0.505
```

The classification procedure which assigns the subject into a group for which the posterior median of the individual component probability is maximal is implemented as follows.

```
R> groupMed <- apply(mod[[1]]$quant.comp.prob[["50%"]], 1, which.max)
R> pMed <- apply(mod[[1]]$quant.comp.prob[["50%"]], 1, max)
R> table(groupMed)
```

```
groupMed
  1    2
162  98
```

```
R> table(groupMean, groupMed)
```

```
      groupMed
groupMean  1    2
  1 161    0
  2   1  98
```

Classification of one subject (id 66) has changed. Nevertheless, both the posterior mean and the posterior median of  $p_{i,g(i)}$  are close to the classification threshold of 0.5 as we illustrate using the following code:

```
R> pMeanMed <- data.frame(Mean = pMean, Median = pMed)
R> rownames(pMeanMed) <- unique(PBC910$id)
R> print(pMeanMed[groupMean != groupMed, ])
```

```
      Mean Median
66 0.502  0.505
```

### *Posterior distribution and credible intervals of the individual component probabilities*

Up to now, only one-number characteristics of the posterior distributions of the ICPs were used for the purposes of clustering which in fact corresponds to the use of the point estimates in a frequentist setting. Nevertheless, the one-number characteristics have a different informative value for different subjects due to different variability of the posterior distributions, i.e., as point estimates they have a different precision when talking in frequentist terms. To illustrate this, Figure 5 shows histograms of sampled values of  $p_{i,1}$  (i.e., their estimated posterior densities) for three subjects (id 2, 7, 11) together with their posterior mean and median. The figure was created by applying the standard R function `hist` to appropriate columns of the

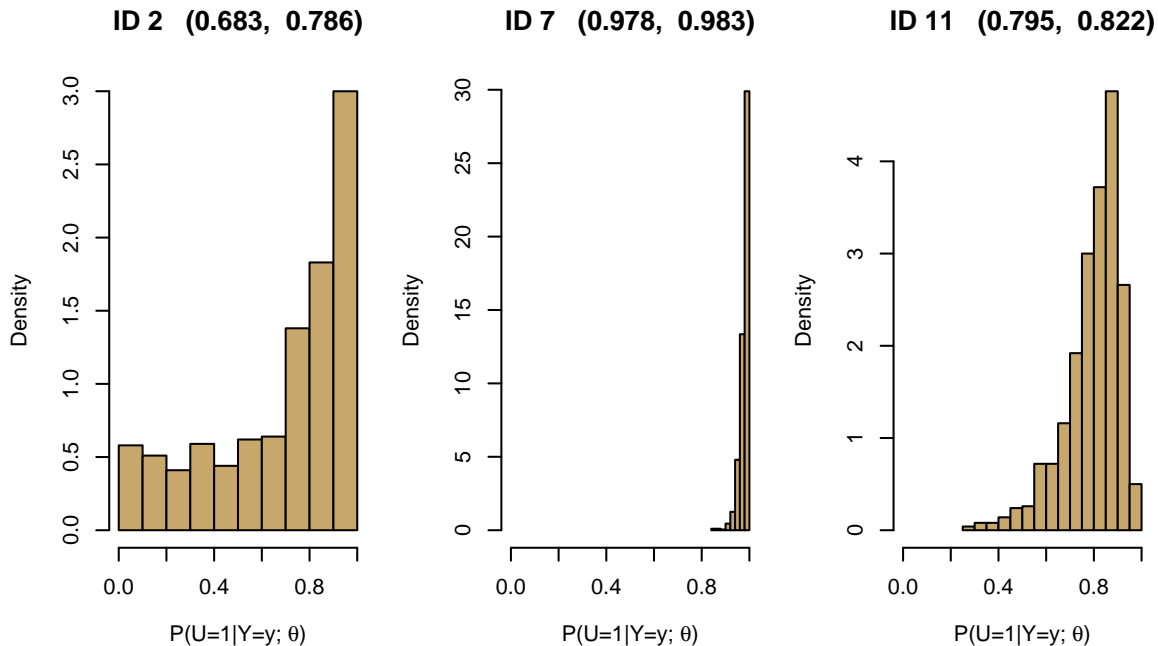


Figure 5: Histograms of sampled values of the individual component probabilities  $p_{i,1}(\theta)$  for three subjects. Above the plot: posterior mean and posterior median of  $p_{i,1}(\theta)$ .

`mod[[1]]$comp.prob` matrix. Besides a different variability of the posterior distributions of the ICPs, Figure 5 also shows that those posterior distributions are often skewed which lead us to conclude that if only a one-number characteristic of the posterior distribution should be used, then the posterior medians are probably a better choice than the posterior means.

When also taking into account the variability of the posterior distribution of the ICPs, we also suggest incorporating the credible intervals of the ICPs in the classification. We ultimately classify the  $i$ th subject into group  $g(i)$  only if  $\hat{p}_{i,g(i)} = \max_{k=1,\dots,K} \hat{p}_{i,k}$  ( $\hat{p}_{i,k}$  are suitable one-number characteristics of the posterior distributions of the ICPs, e.g., the posterior means or medians) and at the same time the lower limit of the credible interval for  $p_{i,g(i)}$  exceeds a certain threshold, for instance 0.5. Indeed, some subjects may remain unclassified with such a rule. Nevertheless, this may not be a problem in practice if we view this classification exercise as one step in a multi-level classification procedure. Subjects with unknown cluster pertinence (and only those subjects) may proceed to another (more complicated or more expensive) level of the full classification procedure. A classification which exploits the HPD credible intervals (calculated using the `cod` function `HPDinterval()`) for the ICPs with a threshold for ultimate classification of 0.5 is implemented by the following code. The missing value indicator `NA` is used to mark the classification for those subjects for whom none of the lower limits of the credible intervals exceeds the threshold.

```
R> pHPD <- HPDinterval(mcmc(mod[[1]]$comp.prob))
R> pHPDlower <- matrix(pHPD[, "lower"], ncol = 2, byrow = TRUE)
R> pHPDupper <- matrix(pHPD[, "upper"], ncol = 2, byrow = TRUE)
R> rownames(pHPDlower) <- rownames(pHPDupper) <- unique(PBC910$id)
R> groupHPD <- groupMed
```

```
R> groupHPD[groupHPD == 1 & pHPDlower[, 1] <= 0.5] <- NA
R> groupHPD[groupHPD == 2 & pHPDlower[, 2] <= 0.5] <- NA
R> table(groupHPD, useNA = "ifany")
```

```
groupHPD
  1    2 <NA>
123  70   67
```

That is, applying the above-described classification procedure based on the credible intervals for the ICPs lead to 123 subjects being classified in group 1 and 70 subjects being classified in group 2. For 67 subjects, the ultimate classification cannot be determined due to the fact that their observed longitudinal markers do not provide enough certainty for classification into any of the two considered groups.

#### *Observed longitudinal profiles by clusters*

Results of the clustering procedure may be visualized by plotting the observed longitudinal profiles using different colors according to the subjects' classification. To achieve this for a classification stored in a vector `groupHPD`, the following procedure leading to Figure 6 can be used. We first add to the original `data.frame` `PBC910` factors created from the classification vectors `groupMed` and `groupHPD` while turning NA into one of the `groupMed` factor levels.

```
R> TAB <- table(PBC910$id)
R> PBC910$groupMed <- factor(rep(groupMed, TAB))
R> PBC910$groupHPD <- factor(rep(groupHPD, TAB), exclude = NULL)
```

Second, we again extract the longitudinal profiles while keeping also the group indicators in the resulting object.

```
R> ip <- getProfiles(t = "month", y = c("lbili", "platelet", "spiders",
+   "jspiders", "groupHPD"), id = "id", data = PBC910)
```

Finally, we plot the observed longitudinal profiles while using different colors for profiles of subjects belonging to different groups according to variable `groupHPD`. To this end, arguments `gvar` and `col` of the `mixAK` function `plotProfiles()` are exploited (only code for upper left panel of Figure 6 with the variable `lbili` shown):

```
R> GCOL <- rainbow_hcl(3, start = 220, end = 40, c = 50, l = 60)[c(2, 3, 1)]
R> names(GCOL) <- levels(PBC910$groupHPD)
R> plotProfiles(ip = ip, data = PBC910, var = "lbili", tvar = "month",
+   gvar = "groupHPD", col = GCOL, auto.layout = FALSE,
+   main = "Log(bilirubin)", xlab = "Time (months)",
+   ylab = "Log(bilirubin)")
```

Figure 6 now shows even better than the previously discussed Figure 4 that group 1 is composed of subjects with lower bilirubin values compared to group 2. Further, we see that the platelet count values do not contribute greatly to the classification and that the occurrence of blood vessel malformations is visibly lower in group 1 than in group 2. All these findings are

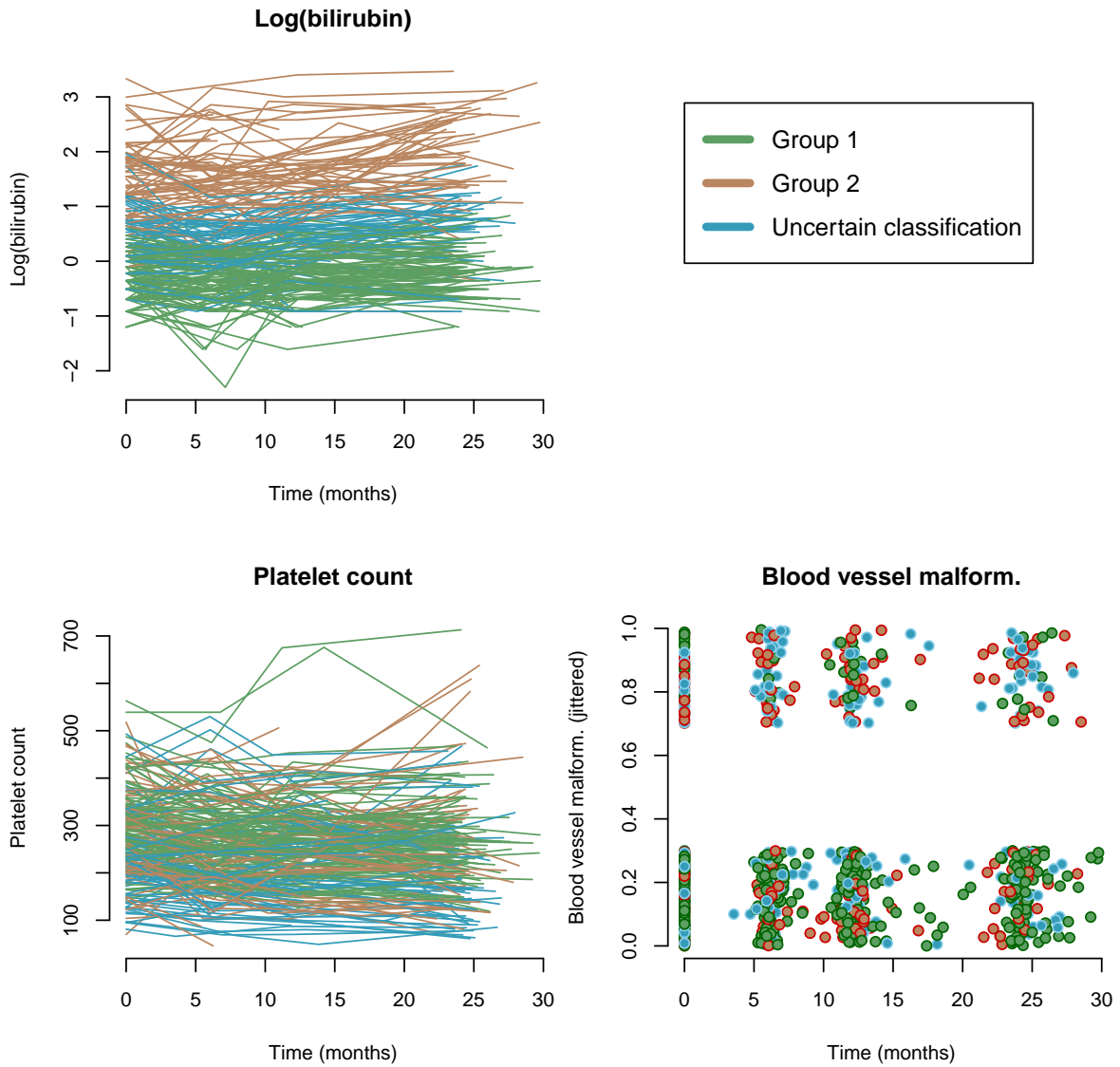


Figure 6: Observed longitudinal profiles of considered markers colored according to classification based on the HPD credible intervals for the individual component probabilities.

indeed in accordance with our findings based on the estimated cluster specific mean profiles discussed earlier for Figure 4. Probably a clearer picture can be obtained by plotting the longitudinal profiles of subjects belonging to different groups in separate plots, the R code of which is available in the supplementary material of the paper.

### 3.8. Selection of a number of mixture components

The final part of our example analysis shall be devoted to the selection of a suitable number of mixture components using the approaches outlined in Section 2.4.

*Penalized expected deviance*

For a fitted model, the value of the penalized expected value (14) together with related quantities is available as a vector element `PED` of an object of class ‘`GLMM_MCMClist`’:

```
R> print(mod$PED)
```

```
D.expect  p(opt)      PED  wp(opt)      wPED
14088.3   74.5  14162.8   74.8  14163.1
```

The same values were also seen on top of the output from `print(mod)` in Section 3.6. In the output above, `D.expect` is the estimated value of  $E\{D(\boldsymbol{\theta}) \mid \mathbf{Y} = \mathbf{y}\}$  from Equation 14, based on both sampled chains stored in the ‘`GLMM_MCMClist`’ object `mod`. Further, `p(opt)` and `wp(opt)` are the values of the penalty term  $p_{opt}$  estimated by the importance sampling as generally outlined by Plummer (2008). Both `p(opt)` and `wp(opt)` are weighted averages of certain quantities based on generated importance samples, where `p(opt)` exploits the unity weights whereas `wp(opt)` the weights proposed by Plummer (2008). Nevertheless, as argued by him, for regular models the difference between `p(opt)` and `wp(opt)` becomes vanishingly small (as it is the case in our example) with increasing sample size. Therefore, in many practical situations it does not really matter which method is used to calculate the penalty term  $p_{opt}$ . Finally, `PED` and `wPED` are the PED values calculated using Equation 14 where `p(opt)` and `wp(opt)`, respectively, are used as the value of  $p_{opt}$ .

*Deviance samples*

In Section 2.4 we further mentioned that the selection of a number of mixture components can also be based on the posterior distribution of the deviances. The two parallel samples from this posterior distribution which we later use to calculate either the posterior probability (15) or the values of the CDF (17), are available as `Deviance1` and `Deviance2` vector elements of an object of class ‘`GLMM_MCMClist`’. For example, the first ten values of the first sample are:

```
R> print(mod$Deviance1[1:10], digits = 9)
```

```
[1] 14086.4445 14095.8681 14113.7117 14097.6682 14109.1799 14102.2900
[7] 14110.1603 14098.9010 14109.3791 14096.8414
```

Note that the `mod$Deviance1` vector provides the same posterior sample as the `mod[[1]]$Deviance` vector (see Section 3.5) and analogously `mod$Deviance2` vector provides the same sample as the `mod[[2]]$Deviance` vector. Nevertheless, when comparing the output above from that on page 15 where we printed the first ten elements of a vector `mod[[1]]$Deviance`, a negligible difference is visible. This arises from the fact that the Laplace approximation around the mode of the integrand located by the Newton-Raphson algorithm is used to calculate numerically the integrals in (7) needed for the deviance evaluation and the initial values for the Newton-Raphson algorithm are different for calculation of `mod[[1]]$Deviance` and `mod$Deviance1`, respectively. Analogously, the values in `mod$Deviance2` are slightly different from those in `mod[[2]]$Deviance`.

*Comparison of models with different numbers of components*

To use the PED, or any other characteristic of the fitted model for the selection of a number of mixture components, models with different values of  $K$  must be first fitted. This can be done by running the following code where we run the posterior MCMC simulation for models with  $K = 1, 2, 3, 4$ . To use only a necessary amount of the operating memory we keep from each model only the deviance samples (to be stored in the lists `Devs1` and `Devs2`) and the values of the penalized expected deviance and the related quantities (to be stored in a `data.frame` `PED`).

```
R> Devs1 <- Devs2 <- list()
R> PED <- data.frame()
R> for (K in 1:4) {
+   set.seed(20042005 + K)
+   modK <- GLMM_MCMC(y = PBC910[, c("lbili", "platelet", "spiders")],
+     dist = c("gaussian", "poisson(log)", "binomial(logit)",
+     id = PBC910[, "id"], x = list(lbili = "empty", platelet = "empty",
+     spiders = PBC910[, "month"]), z = list(lbili = PBC910[, "month"],
+     platelet = PBC910[, "month"], spiders = "empty"),
+     random.intercept = rep(TRUE, 3), prior.b = list(Kmax = K),
+     nMCMC = c(burn = 100, keep = 1000, thin = 10, info = 100),
+     parallel = FALSE)
+   Devs1[[K]] <- modK$Deviance1
+   Devs2[[K]] <- modK$Deviance2
+   PED <- rbind(PED, modK$PED)
+   colnames(PED) <- names(modK$PED)
+ }
```

Consequently, the PED values for the four models are as follows, leading us to conclude that the best model having the lowest value of the PED is that with  $K = 2$  mixture components:

```
R> print(PED, digits = 6)
```

	D.expect	p(opt)	PED	wp(opt)	wPED
1	14242.0	36.0238	14278.0	35.8741	14277.8
2	14088.3	74.5441	14162.8	74.8084	14163.1
3	14057.3	131.5205	14188.8	132.7313	14190.0
4	14032.7	179.1014	14211.8	277.4600	14310.2

The values of the posterior probabilities (15) and (16) recommended for the model comparison by Aitkin *et al.* (2009) are easily obtained by using another **mixAK** function `summaryDiff()` applied to the posterior samples of the model deviances. In the following, let  $D_K(\theta)$  denote the deviance of a model with  $K$  mixture components. We first compare models with  $K = 2$  and  $K = 1$ :

```
R> sD21 <- summaryDiff(c(Devs1[[2]], Devs2[[2]]), c(Devs1[[1]], Devs2[[1]]))
R> print(sD21, digits = 4)
```



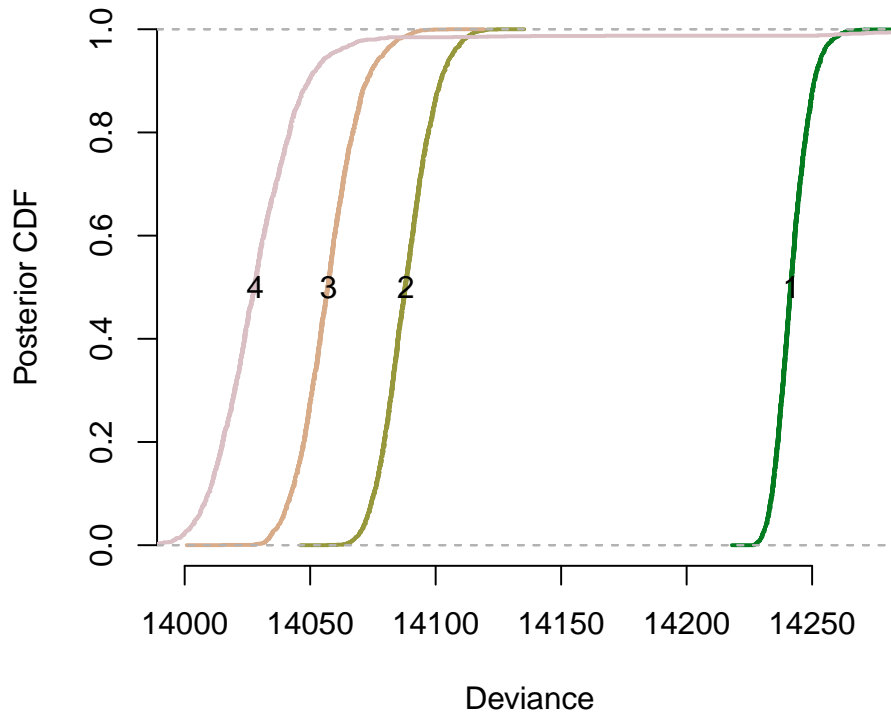


Figure 7: Posterior cumulative distribution functions of the observed data deviances for models with  $K = 1, 2, 3, 4$ .

```
$summary
```

```
  Mean  2.5%   50%  97.5%
-153.7 -177.4 -154.2 -129.4
```

```
$Pcut
```

```
P(diff < -4.39)   P(diff < 0)
                1             1
```

The output first provides posterior summary statistics for the difference between  $D_2(\boldsymbol{\theta}) - D_1(\boldsymbol{\theta})$  and then the values of the posterior probabilities (16) and (15), respectively. As both of these probabilities are practically equal to one, it is clear that a model with  $K = 2$  is better than that with  $K = 1$  and hence there is a strong evidence that at least two clusters are present in the data at hand. Using a simple loop (code not shown), analogous values for a comparison of each pair of models were calculated:

	P(diff < -4.39)	P(diff < 0)	Mean	2.5%	50%	97.5%
2 - 1	1.000	1.000	-153.7	-177.4	-154.2	-129.401
3 - 1	1.000	1.000	-184.7	-210.6	-184.8	-155.038
4 - 1	0.988	0.988	-209.3	-245.7	-214.1	-170.617
3 - 2	0.952	0.976	-31.0	-62.7	-31.1	-0.121
4 - 2	0.981	0.982	-55.6	-95.5	-60.3	-13.360
4 - 3	0.884	0.914	-24.6	-66.1	-28.9	16.531

If the approach of [Aitkin \*et al.\* \(2009\)](#) is used for model selection, then we would choose

a model with  $K = 3$  components. The posterior probability (16) shown in the first column of the output above is clearly higher than the value of 0.9 advocated by Aitkin *et al.* (2009) to declare model with  $K = 3$  components being evidently better than a model with  $K = 2$  components. Nevertheless, the value of (16) drops below 0.9 when comparing  $K = 4$  to  $K = 3$  components.

Finally, the graphical comparison of the CDFs of the deviances of the competing models suggested by Aitkin (2010) is provided by Figure 7 prepared using the following code:

```
R> COL <- terrain_hcl(4, c = c(65, 15), l = c(45, 80), power = c(0.5, 1.5))
R> plot(c(14000, 14275), c(0, 1), type = "n", xlab = "Deviance",
+   ylab = "Posterior CDF")
R> for (K in 1:4) {
+   medDEV <- median(c(Devs1[[K]], Devs2[[K]]))
+   ECDF <- ecdf(c(Devs1[[K]], Devs2[[K]]))
+   plot(ECDF, col = COL[K], lwd = 2, add = TRUE)
+   text(medDEV + 0.5, 0.5, labels = K)
+ }
```

Figure 7 shows a huge improvement of the model with respect to its deviance when moving from the  $K = 1$  model to the  $K = 2$  model. The variability of the posterior distribution of the deviance in a model with  $K = 2$  components is practically the same as with  $K = 1$ . Nevertheless, the  $K = 2$  deviance posterior distribution is clearly shifted to left compared to  $K = 1$ . Almost the same conclusion can be drawn when comparing models with  $K = 3$  and  $K = 2$  components, and also  $K = 4$  and  $K = 3$  components.

## 4. Conclusions and outlook

In particular, on a practical example, this paper provides an overview of the capabilities of the **mixAK** package suitable for clustering based on multivariate continuous and discrete longitudinal data. It was not possible to describe here all the options and arguments of the respective functions. Nevertheless, a detailed description of them is available in the help pages of the package. Additional technical details can also be found in the package vignettes.

Further, this manuscript presents the **mixAK** package mainly as a clustering tool. Nevertheless, the package can also be used in situations when purely a regression analysis with multivariate longitudinal outcomes is of interest. In this case, Equation 8 suggests that the (multivariate) mixture GLMM can be considered as a version of the (multivariate) GLMM being robustified against misspecification of the random effects distribution by assuming a normal mixture rather than the conventionally used normal distribution.

The functionality of the **mixAK** package can be extended in the future in several other directions to provide even more flexible clustering tools. First, we may restrict the component covariance matrices  $\mathbb{D}_1, \dots, \mathbb{D}_K$  in (4) to be common as it is often assumed in mixture applications. Other restrictions imposed on the component covariance matrices, e.g., those considered by the R package **HDclassif**, can be thought of as well, of course. This would allow the user to consider a more parsimonious model underlying the clustering procedure. Second, additional flexibility of the clustering procedure might be achieved by assuming that not only the distributional components of the random effects vectors of the underlying generalized linear

mixed model are cluster specific but that also the fixed effects vectors  $\alpha_1, \dots, \alpha_R$  vary across clusters. Further, a normal distribution in (4) could be replaced by, e.g., the multivariate  $t$  distribution or even more flexible multivariate skew  $t$  distribution (Azzalini and Capitanio 2003) to have a model being able to capture outlying observations.

In this manuscript, we have shown an analysis with outcomes of three types (dichotomous, count, and continuous) for which three benchmark GLMM's, mentioned at the end of Section 2.1, were assumed. Currently, those are the only GLMM's implemented by the **mixAK** package. One important type of outcome, namely a genuine categorical (multinomial) response, is thus not yet covered by the package. A possible step towards the ability to include also the multinomial response could consist of incorporating a mixed-effects version of any of routinely used regression models for such type of outcome in assumption (A1), e.g., the baseline-category logit model (Agresti 2002, Chapter 7). Nevertheless, this goes far beyond the scope of this paper.

Finally, it is worth mentioning that the methods implemented in the **mixAK** package were primarily developed having classical (bio)statistical applications in mind. Notably, even though multivariate outcomes ( $R > 1$ ) are considered, the presented methods and their **mixAK** implementation, if not properly adjusted, are likely not suitable for the analysis of ultra-dimensional outcomes ( $R$  being huge) that are often encountered, e.g., in the context of bioinformatics.

## Computational details

The output shown in this article was obtained using R version 3.1.1 (2014-07-10), **mixAK** 3.8, and the following contributed packages which are the dependencies or imports of the package **mixAK**: **colorspace** 1.2-4 (Ihaka, Murrell, Hornik, Fisher, and Zeileis 2013; Zeileis, Hornik, and Murrell 2009), **lme4** 1.1-7 (Bates, Maechler, Bolker, and Walker 2014), **fastGHQuad** 0.1-1 (Blocker 2011), **mnormt** 1.5-1 (Azzalini and Genz 2014). Finally, routines from **coda** 0.16-1 (Plummer *et al.* 2006) were used in this manuscript.

## Acknowledgments

We would like to thank Mr. Steven Riley for his help with English corrections of the manuscript. Last but not the least, we thank two anonymous referees and an Associate Editor for very detailed and stimulative comments to the earlier versions of this manuscript that led to its considerable improvement.

Financial support of the first author from the Interuniversity Attraction Poles Programme (IAP-network P7/06), Belgian Science Policy Office, is gratefully acknowledged. The second author was supported by the Czech Science Foundation grant GAČR P403/12/1557.

This document was prepared using **Sweave** (Leisch 2002).

## References

Agresti A (2002). *Categorical Data Analysis*. 2nd edition. John Wiley & Sons, Hoboken.

- Aitkin M (2010). *Statistical Inference: An Integrated Bayesian/Likelihood Approach*. Chapman & Hall/CRC, Boca Raton.
- Aitkin M, Liu CC, Chadwick T (2009). “Bayesian Model Comparison and Model Averaging for Small-Area Estimation.” *The Annals of Applied Statistics*, **3**(1), 199–221.
- Andrews JL, McNicholas PD (2012). “Model-Based Clustering, Classification, and Discriminant Analysis via Mixtures of Multivariate  $t$ -Distributions.” *Statistics and Computing*, **22**(5), 1021–1029.
- Andrews JL, McNicholas PD (2013). *teigen: Model-Based Clustering and Classification with the Multivariate  $t$ -Distribution*. R package version 2, URL <http://CRAN.R-project.org/package=teigen>.
- Auder B, Lebre R, Iovleff S, Langrognet F (2014). *Rmixmod: An Interface for MIXMOD*. R package version 2.0.2, URL <http://CRAN.R-project.org/package=Rmixmod>.
- Azzalini A, Capitanio A (2003). “Distributions Generated by Perturbation of Symmetry with Emphasis on a Multivariate Skew  $t$ -Distribution.” *Journal of the Royal Statistical Society B*, **65**(2), 367–389.
- Azzalini A, Genz A (2014). *mnormt: The Multivariate Normal and  $t$  Distributions*. R package version 1.5-1, URL <http://azzalini.stat.unipd.it/SW/Pkg-mnormt>.
- Bates D, Maechler M, Bolker B, Walker S (2014). *lme4: Linear Mixed-Effects Models using Eigen and S4*. R package version 1.1-7, URL <http://CRAN.R-project.org/package=lme4>.
- Benaglia T, Chauveau D, Hunter DR, Young D (2009). “**mixtools**: An R Package for Analyzing Finite Mixture Models.” *Journal of Statistical Software*, **32**(6), 1–29. URL <http://www.jstatsoft.org/v32/i06/>.
- Bergé L, Bouveyron C, Girard S (2012). “**HDclassif**: An R Package for Model-Based Clustering and Discriminant Analysis of High-Dimensional Data.” *Journal of Statistical Software*, **46**(6), 1–29. URL <http://www.jstatsoft.org/v46/i06/>.
- Blocker AW (2011). *fastGHQuad: Fast Rcpp Implementation of Gauss-Hermite Quadrature*. R package version 0.1-1, URL <http://CRAN.R-project.org/package=fastGHQuad>.
- Bock HH (1996). “Probabilistic Models in Cluster Analysis.” *Computational Statistics & Data Analysis*, **23**(1), 5–28.
- Celeux G, Forbes F, Robert CP, Titterton DM (2006). “Deviance Information Criteria for Missing Data Models.” *Bayesian Analysis*, **1**(4), 651–706.
- Dickson ER, Grambsch PM, Fleming TR, Fisher LD, Langworthy A (1989). “Prognosis in Primary Biliary-Cirrhosis: Model for Decision-Making.” *Hepatology*, **10**(1), 1–7.
- Fraley C, Raftery AE (2002). “Model-Based Clustering, Discriminant Analysis, and Density Estimation.” *Journal of the American Statistical Association*, **97**(458), 611–631.

- Fraley C, Raftery AE, Murphy TB, Scrucca L (2012). “**mclust** Version 4 for R: Normal Mixture Modeling for Model-Based Clustering, Classification, and Density Estimation Technical Report No. 597.” *Technical report*, Department of Statistics, University of Washington.
- Gelfand AE, Sahu SK, Carlin BP (1995). “Efficient Parametrisations for Normal Linear Mixed Models.” *Biometrika*, **82**(3), 479–499.
- Grün B, Leisch F (2008). “FlexMix Version 2: Finite Mixtures with Concomitant Variables and Varying and Constant Parameters.” *Journal of Statistical Software*, **28**(4), 1–35. URL <http://www.jstatsoft.org/v28/i04>.
- Ihaka R, Murrell P, Hornik K, Fisher JC, Zeileis A (2013). *colorspace: Color Space Manipulation*. R package version 1.2-4, URL <http://CRAN.R-project.org/package=colorspace>.
- Komárek A (2009). “A New R Package for Bayesian Estimation of Multivariate Normal Mixtures Allowing for Selection of the Number of Components and Interval-Censored Data.” *Computational Statistics & Data Analysis*, **53**(12), 3932–3947.
- Komárek A (2014). *mixAK: Multivariate Normal Mixture Models and Mixtures of Generalized Linear Mixed Models Including Model Based Clustering*. R package version 3.8, URL <http://CRAN.R-project.org/package=mixAK>.
- Komárek A, Hansen BE, Kuiper EMM, van Buuren HR, Lesaffre E (2010). “Discriminant Analysis Using a Multivariate Linear Mixed Model with a Normal Mixture in the Random Effects Distribution.” *Statistics in Medicine*, **29**(30), 3267–3283.
- Komárek A, Komárková L (2013). “Clustering for Multivariate Continuous and Discrete Longitudinal Data.” *The Annals of Applied Statistics*, **7**(1), 177–200.
- Leisch F (2002). “Dynamic Generation of Statistical Reports Using Literate Data Analysis.” In W Härdle, B Rönz (eds.), *COMPSTAT 2002 – Proceedings in Computational Statistics*, pp. 575–580. Physica-Verlag, Heidelberg.
- Macdonald P, Du J (2012). *mixdist: Finite Mixture Distribution Models*. R package version 0.5-4, URL <http://CRAN.R-project.org/package=mixdist>.
- McNicholas PD, Jampani KR, Subedi S (2012). *longclust: Model-Based Clustering and Classification for Longitudinal Data*. R package version 1.1, URL <http://CRAN.R-project.org/package=longclust>.
- Plummer M (2008). “Penalized Loss Functions for Bayesian Model Comparison.” *Biostatistics*, **9**(3), 523–539.
- Plummer M, Best N, Cowles K, Vines K (2006). “**coda**: Convergence Diagnosis and Output Analysis for MCMC.” *R News*, **6**(1), 7–11. URL <http://CRAN.R-project.org/doc/Rnews/>.
- Proust-Lima C, Philipps V, Diakite A, Lique B (2014). *lcmm: Estimation of Latent Class Mixed Models, Joint Latent Class Mixed Models and Mixed Models for Curvilinear Outcomes*. R package version 1.6.4, URL <http://CRAN.R-project.org/package=lcmm>.

- R Core Team (2014). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <http://www.R-project.org/>.
- Rubin DB (1976). “Inference and Missing Data.” *Biometrika*, **63**(3), 581–592.
- Spiegelhalter DJ, Best NG, Carlin BP, van der Linde A (2002). “Bayesian Measures of Model Complexity and Fit.” *Journal of the Royal Statistical Society B*, **64**(4), 583–639.
- Stephens M (2000). “Dealing with Label Switching in Mixture Models.” *Journal of the Royal Statistical Society B*, **62**(4), 795–809.
- Zeileis A, Hornik K, Murrell P (2009). “Escaping RGBland: Selecting Colors for Statistical Graphics.” *Computational Statistics & Data Analysis*, **53**(9), 3259–3270.
- Zhong W, Ma P (2012). *MFDA: Model Based Functional Data Analysis*. R package version 1.1-4, URL <http://CRAN.R-project.org/package=MFDA>.

**Affiliation:**

Arnošt Komárek  
Department of Probability and Mathematical Statistics  
Faculty of Mathematics and Physics  
Charles University in Prague  
Sokolovská 83  
CZ–186 75 Praha 8 – Karlín, Czech Republic  
E-mail: [arnost.komarek@mff.cuni.cz](mailto:arnost.komarek@mff.cuni.cz)  
URL: <http://msekce.karlin.mff.cuni.cz/~komarek/>

Lenka Komárková  
Department of Exact Methods  
Faculty of Management  
University of Economics in Prague  
Jarošovská 1117  
CZ–377 01 Jindřichův Hradec, Czech Republic  
E-mail: [komarkol@fm.vse.cz](mailto:komarkol@fm.vse.cz)