

Classification Based on Multivariate Mixed Type Longitudinal Data with an application to the EU-SILC database

Jan Vávra · Arnošt Komárek

Received: June 1, 2021 / Accepted: Month DD, YYYY

Abstract Although, many nowadays studies gather data of diverse nature (numeric quantities, binary indicators or ordered categories) on the same units repeatedly over time, there exist only limited number of approaches in the literature to analyse so called *mixed-type* longitudinal data. We present a statistical model capable of joint modelling of several mixed-type outcomes, which also accounts for possible dependencies among investigated outcomes. A thresholding approach to link binary or ordinal variables to their latent numeric counterparts allows us to jointly model all, including latent, numeric outcomes using a multivariate version of the linear mixed-effects model. We avoid the independence assumption over outcomes by relaxing the variance matrix of random effects to a completely general positive definite matrix. Moreover, we follow Model Based Clustering (MBC) methodology to create a mixture of such models to model heterogeneity in temporal evolution of considered outcomes. The estimation of such hierarchical model is approached by Bayesian principles with the use of Markov Chain Monte Carlo (MCMC) methods. After a successful simulation study with the aim to examine the ability to consistently estimate the true parameter values and thus discovering the different patterns, the EU-SILC dataset consisting of Czech households each followed for four years in a time span 2005–2016 was analysed. Households were classified into groups with similar evolution of several closely related indicators of monetary poverty based on estimated classification probabilities.

Keywords Multivariate longitudinal data · Mixed type outcome · Model based clustering · Classification · EU-SILC

1 Introduction

In different types of studies data are nowadays routinely gathered repeatedly over time on the same units leading to *longitudinal* or *panel* data. On top of that, multiple outcomes, both *numeric* and *categorical*, i.e., of a *mixed type*, are recorded at each measurement occasion leading to *multivariate mixed type longitudinal data*. An example of such a dataset which also motivates our research is *The European Union Statistics on Income and Living Conditions database* (EU-SILC, <https://ec.europa.eu/eurostat/web/microdata/european-union-statistics-on-income-and-living-conditions>). This is an instrument with the goal to collect timely and comparable cross-sectional and longitudinal multidimensional microdata on income, poverty, social exclusion and living conditions in the European Union, Iceland, Norway and Switzerland. The reference population includes all private households of respective countries and variables, which are collected annually via questionnaires, refer both to households and to individuals from the household. In this paper, we concentrate on household specific data from the Czech Republic (period 2005–2016), each household was followed annually for a period of 4 years. In total,

Jan Vávra
Faculty of Mathematics and Physics, Charles University, Prague, Czech Republic
Tel.: +420 951 553 385
E-mail: vavraj@karlin.mff.cuni.cz

Arnošt Komárek
Faculty of Mathematics and Physics, Charles University, Prague, Czech Republic
Tel.: +420 951 553 282
E-mail: komarek@karlin.mff.cuni.cz

$n = 20323$ households will be analyzed. The aim of the research is to find typical patterns of temporal evolution of several indicators related to poverty and material deprivation. The relevant outcome variables are not only *numeric* (e.g., income) but also *binary* (e.g., ability to cover unexpected expenses), or *ordinal* (e.g., level of a financial burden of housing). From a data analytic point of view, it is our aim to develop a clustering approach suitable for longitudinal data of a mixed type that allows for above mentioned types of outcome variables.

To formalize the task, we are assuming that data are composed of n independently behaving units (e.g., households) and for the i th unit ($i = 1, \dots, n$), in total R outcome variables $Y_{i,j}^r$ ($r = 1, \dots, R, j = 1, \dots, n_i$) are being gathered at each of n_i measurement occasions that take place at times $t_{i,1}, \dots, t_{i,n_i}$. On top of that, each outcome variable $Y_{i,j}^r$ might be either *numeric*, *binary* or *ordinal*, see Figure 1 for an example. Finally, each observation might be supplemented by a vector $\mathbf{v}_{i,j}^r$ of additional covariates that may explain the outcome variability. In summary, the i th unit is represented by data $\mathcal{D}_i = \{Y_{i,j}^r, \mathbf{v}_{i,j}^r, t_{i,j} : r = 1, \dots, R, j = 1, \dots, n_i\}, i = 1, \dots, n$ and the task is to use this information to classify each unit into one of $K > 1$ groups with a priori unknown structure.

Due to complexity of a data structure and also due to the fact that possibly different numbers n_i of measurement occasions appear in data for different units, classical distance-based clustering methods like hierarchical clustering or the K -means method and their many extensions (see, e.g., [Hastie et al, 2009](#), Chapters 13 and 14) could hardly be used. On the other hand, methods that further develop ideas of model based clustering (MBC, [Banfield and Raftery, 1993](#); [Fraley and Raftery, 2002](#)) and that exploit mixtures of suitable statistical models proved to be useful in similar situations. A classical model to analyze *continuous* longitudinal outcomes is the linear mixed model (LMM, [Laird and Ware, 1982](#)) and hence not surprisingly, several MBC procedures based on mixtures of LMM's appeared in the literature. A work by [Verbeke and Lesaffre \(1996\)](#), where growth curves are classified, provides one of the first methods of this type even though not explicitly called MBC at that time. More recently, an application of similar ideas to clustering of gene-expression data is covered by [Celeux et al \(2005\)](#). Subsequently, [De la Cruz-Mesía et al \(2008\)](#) base their MBC procedure for longitudinal data on a non-linear mixed model. The situation of more than one ($R > 1$) outcome being available for the clustering, nevertheless, all of them still being *continuous*, is considered by [Villarroel et al \(2009\)](#).

The MBC methods developed for functional data and (continuous) stochastic processes could also be employed if we keep dealing with *continuous* and moreover univariate ($R = 1$) longitudinal data (e.g. [James and Sugar, 2003](#); [Ma et al, 2006](#); [Liu and Yang, 2009](#); [McNicholas and Murphy, 2010](#)). [Frühwirth-Schnatter \(2011\)](#) provides a comprehensive overview. A possibility to develop the MBC for non-continuous longitudinal data is to replace LMM by a generalized linear mixed model (GLMM) in the underlying mixture of models. See, e.g., [Molenberghs](#)

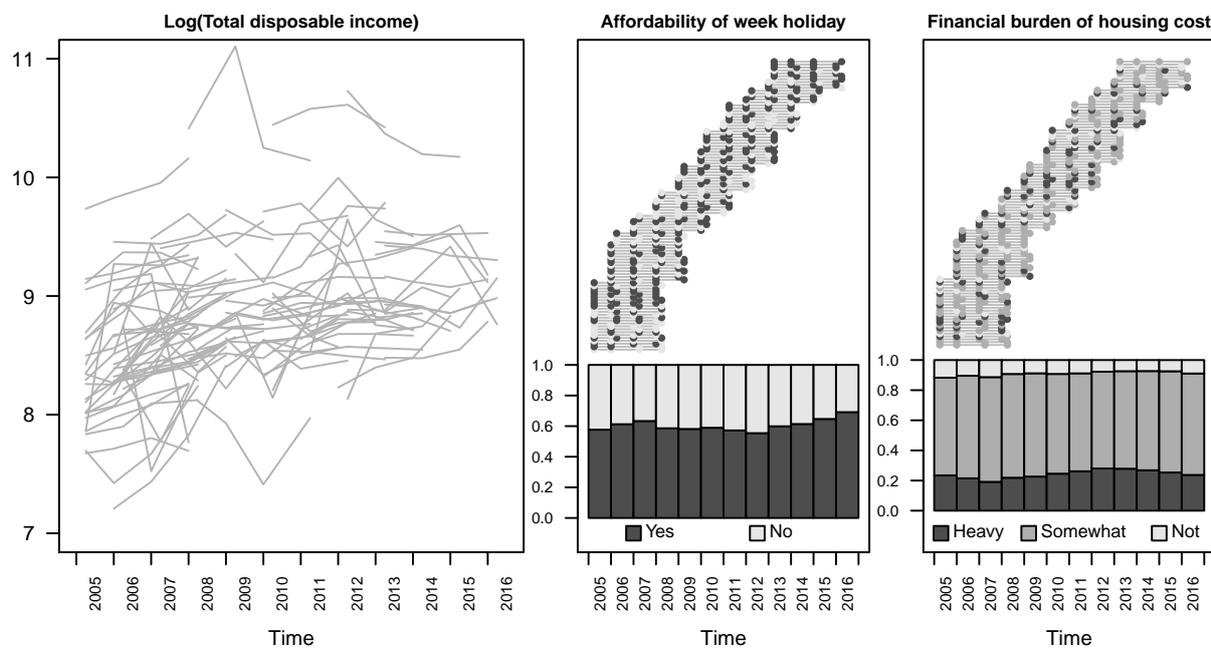


Fig. 1: EU-SILC data (Czech Republic). Observed longitudinal profiles of three (numeric, binary and ordinal) outcomes.

and Verbeke (2005, Chapter 14) who also provide an example of such a clustering procedure in their Section 23.3. Nevertheless, it is still only possible to use a single ($R = 1$) longitudinal outcome.

On the other hand, only little previous work seems to be available in the literature in case where units are to be classified based on multivariate ($R > 1$) and possibly non-continuous longitudinal data. If all outcomes are of the same type (e.g., all binary), a method based again on a mixture of mixed models is offered by the \mathbb{R} package `lcmm` (Proust-Lima et al, 2017). Nevertheless, for the MBC based on multivariate ($R > 1$) mixed type longitudinal data, the only two approaches we are aware of and that into some extent allow for classification, are those implemented in the \mathbb{R} packages `flexmix` (Grün and Leisch, 2008) and `mixAK` (Komárek and Komárková, 2013, 2014). Nevertheless, both of the two approaches are lacking some important aspects. First, Grün and Leisch (2008) assume independence of different longitudinal outcomes measured at one occasion. This may not only be unrealistic but also prevents the analyst from exploiting information provided by the dependence structure among the R outcomes in the clustering procedure. Even though certain form of dependence is considered by Komárek and Komárková (2013, 2014), only binary or count non-continuous outcomes are considered which does not allow for use with a typical questionnaire data like the EU-SILC database where many outcome variables are of an ordinal nature.

One of the reasons why not much is available to perform clustering based on multivariate mixed type longitudinal data is perhaps also the fact that even statistical models needed to develop the MBC procedure that would allow for datasets of a considered structure are relatively scarce in literature. This especially if we seek for models that realistically account for possible dependencies between different outcome variables gathered at one occasion. Fieuws and Verbeke (2004) covered in detail a bivariate case of longitudinal data and also in this manuscript, we follow their suggestion to use a multivariate mixed model while specifying a general covariance matrix for the joint distribution of all involved random effects. Later, Fieuws and Verbeke (2006) extended this approach to more than two outcomes by pairwise fitting and construction of pseudo-likelihood to avoid computational problems with covariance matrix of a high dimension. Nevertheless, MBC was not employed in any of those solutions. Recently, Bruckers et al (2016) invented a clustering algorithm that updates pseudo-log-likelihood of the pairwise approach and reclassifies individuals until no change is made. This solution, however, lacks inclusion of binary and ordinal outcomes that we aim to provide in this article.

The rest of the paper is organized as follows. In Section 2 we first outline the approach being capable of a joint modelling of mixed-type (numeric, binary and ordinal) longitudinal data. Second, in Section 3, we incorporate the developed model within the clustering procedure that allows usage of data with a structure analogous to that on Figure 1 and classification of study units into groups with apriori unknown structure. Yet, Section 3 only provides a theoretical clustering concept which assumes a full knowledge of unknown parameters. Transition into a practically applicable procedure is provided in Section 4 which outlines details of a Bayesian approach towards this goal. Further, Section 5 evaluates clustering capabilities of our approach on a simulation study. In Section 6, we apply our method to the EU-SILC database in order to discover clusters of different evolution patterns and to classifying each household. Finally, Section 7 summarizes the proposed model and suggests ways of improvement in reaction to our findings from application.

2 Joint modelling of mixed-type longitudinal data

At each measurement occasion, R outcomes (numeric, ordinal or binary) are observed on each study unit. Let $\mathcal{R} = \{1, \dots, R\} = \mathcal{R}^{\text{Num}} \cup \mathcal{R}^{\text{OB}}$, $\mathcal{R}^{\text{OB}} = \mathcal{R}^{\text{Ord}} \cup \mathcal{R}^{\text{Bin}}$, denote the index set of observed outcomes that consists of indices of numeric outcomes (\mathcal{R}^{Num}), ordinal outcomes (\mathcal{R}^{Ord}) and binary outcomes (\mathcal{R}^{Bin}). Let $\mathbf{Y}_i^r = (Y_{i,1}^r, \dots, Y_{i,n_i}^r)^\top$ be the vector of values of outcome $r \in \mathcal{R}$ of subject $i = 1, \dots, n$ observed at times $\mathbf{t}_i = (t_{i,1}, \dots, t_{i,n_i})$ together with additional covariates $\mathbf{v}_{i,1}^r, \dots, \mathbf{v}_{i,n_i}^r$. Further, let $\mathcal{C}_i^r = \{\mathbf{t}_i, \mathbf{v}_{i,1}^r, \dots, \mathbf{v}_{i,n_i}^r\}$ denote both the measurement times and the covariate values for the outcome r of the i th subject. Finally, let

$$\mathbb{Y}_i = (\mathbf{Y}_i^r, r \in \mathcal{R}), \quad \mathcal{C}_i = \{\mathcal{C}_i^r, r \in \mathcal{R}\} \quad (1)$$

denote all information (outcomes and covariate values) available for the i th subject, which is assumed to be independent of other subjects. \mathbf{Y}^r and \mathcal{C}^r stand for information (outcome and covariate values) regarding one chosen outcome $r \in \mathcal{R}$ from all subjects, while \mathbb{Y} and \mathcal{C} stand for all gathered information (all outcomes and covariate values) from all subjects.

The joint model for data (1) is built hierarchically. It exploits the linear mixed model (LMM) for each longitudinal outcome (each $r \in \mathcal{R}$). In case of binary or ordinal outcomes, the LMM is assumed only latently. Dependencies between different outcomes gathered on a single study unit are captured by considering a vector of shared random effects. In particular, the model is built as follows.

2.1 Numeric longitudinal outcomes

For each numeric outcome $r \in \mathcal{R}^{\text{Num}}$ we directly assume the linear mixed model:

$$\mathbf{Y}_i^r \mid \mathcal{C}_i^r, \mathbf{b}_i^r \sim N_{n_i}(\boldsymbol{\eta}_i^r, \tau_r^{-1} \mathbb{I}_{n_i}), \quad (2)$$

where $\boldsymbol{\eta}_i^r = \mathbb{X}_i^r \boldsymbol{\beta}^r + \mathbb{Z}_i^r \mathbf{b}_i^r$ is the linear predictor consisting of fixed and random effects parts, $\tau_r > 0$ is the precision (inverse variance) of model errors, $\boldsymbol{\beta}^r \in \mathbb{R}^{d^F}$ are fixed effects and $\mathbf{b}_i^r \in \mathbb{R}^{d^R}$ are random effects belonging to subject i . Further, $\mathbb{X}_i^r = (\mathbf{x}_{i,1}^r, \dots, \mathbf{x}_{i,n_i}^r)^\top$ and $\mathbb{Z}_i^r = (\mathbf{z}_{i,1}^r, \dots, \mathbf{z}_{i,n_i}^r)^\top$ are matrices of regressors being derived from the explanatory variables information \mathcal{C}_i^r . For identifiability purposes, matrices \mathbb{X}_i^r and \mathbb{Z}_i^r are assumed not to share the same columns, i.e. created regressor falls exclusively either into fixed effects part or into random effects part of the model.

2.2 Ordinal and binary longitudinal outcomes

The r th binary or ordinal outcome $r \in \mathcal{R}^{\text{OB}}$ is assumed to attain values $0, \dots, L^r - 1$ which are linked to a linear mixed model through the thresholding concept (see, e.g., [Albert and Chib, 1993](#)):

$$Y_{i,j}^r = l \iff \gamma_l^r < Y_{i,j}^{*,r} \leq \gamma_{l+1}^r, \quad (3)$$

where $-\infty = \gamma_0^r < \gamma_1^r < \dots < \gamma_{L^r}^r = \infty$ are (unknown) thresholds categorizing a latent (unobserved) numeric variables $Y_{i,j}^{*,r}$. Let $\boldsymbol{\gamma}^r = (\gamma_1^r, \dots, \gamma_{L^r-1}^r)$ be the vector of thresholds. For identifiability purposes, we will fix one of the thresholds, e.g., the first one γ_1^r while estimating the model. That is, in case of a binary outcome, all threshold parameters are fixed.

Analogously to the case of numeric outcomes, latent numeric variables $Y_{i,j}^{*,r}$ are assumed to follow the linear mixed model

$$\mathbf{Y}_i^{*,r} \mid \mathcal{C}_i^r, \mathbf{b}_i^r \sim N_{n_i}(\boldsymbol{\eta}_i^r, \mathbb{I}_{n_i}) \quad (4)$$

with analogous notation to that in (2). Nevertheless, this time, the precision parameter τ_r of model errors is fixed and equal to one for identifiability purposes.

2.3 Joint model

Let $\mathbb{Y}_i^{\text{N}} = (\mathbf{Y}_i^r, r \in \mathcal{R}^{\text{Num}})$ denote a vector of all numeric outcomes of subject i . Further, let $\mathbb{Y}_i^{*,\text{OB}} = (\mathbf{Y}_i^{*,r}, r \in \mathcal{R}^{\text{OB}})$ be a vector of latent numeric variables being behind all ordinal and binary outcomes. The subvectors of $\mathbb{Y}_i^* := (\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{*,\text{OB}})$ are assumed to follow linear mixed models (2) and (4) with a set of fixed effects $\boldsymbol{\beta} = \{\boldsymbol{\beta}^r, r \in \mathcal{R}\}$ and an overall vector of random effects $\mathbf{b}_i = \{\mathbf{b}_i^r, r \in \mathcal{R}\}$.

In the following, let $\mathbf{b}_i^{\text{N}} = \{\mathbf{b}_i^r, r \in \mathcal{R}^{\text{Num}}\}$ and $\mathbf{b}_i^{\text{OB}} = \{\mathbf{b}_i^r, r \in \mathcal{R}^{\text{OB}}\}$ be random effects related to models for numeric and ordinal/binary longitudinal outcomes, respectively. The overall random effects vector $\mathbf{b}_i \equiv (\mathbf{b}_i^{\text{N}}, \mathbf{b}_i^{\text{OB}})$ is now assumed to follow a multivariate normal distribution with a general covariance matrix, i.e., it is assumed

$$\mathbf{b}_i = \begin{pmatrix} \mathbf{b}_i^{\text{N}} \\ \mathbf{b}_i^{\text{OB}} \end{pmatrix} \stackrel{\text{iid}}{\sim} N_{d^{\text{R}}} \left(\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}^{\text{N}} \\ \boldsymbol{\mu}^{\text{OB}} \end{pmatrix}, \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}^{\text{N}} & \boldsymbol{\Sigma}^{\text{NOB}} \\ \boldsymbol{\Sigma}^{\text{OBN}} & \boldsymbol{\Sigma}^{\text{OB}} \end{pmatrix} \right), \quad (5)$$

where $d^{\text{R}} = d_{\text{N}}^{\text{R}} + d_{\text{OB}}^{\text{R}} = \sum_{r \in \mathcal{R}} d_r^{\text{R}}$ is the total dimension of \mathbf{b}_i , $\boldsymbol{\mu} \in \mathbb{R}^{d^{\text{R}}}$ is the (unknown) mean value of the random effects and $\boldsymbol{\Sigma} > 0$ is the unknown random effects covariance matrix. This matrix is left to be completely general which captures possible dependencies between different longitudinal outcomes.

Throughout the manuscript, the notation $p(\cdot \mid \cdot)$ will stand for a conditional probability distribution function. Next to the fixed effects $\boldsymbol{\beta}$, mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$, the unknown parameters of the model are

$\boldsymbol{\tau} := (\tau_r, r \in \mathcal{R}^{\text{Num}})$, precisions of the error terms of the LMM's for numeric outcomes and $\boldsymbol{\gamma} = \{\gamma^r, r \in \mathcal{R}^{\text{Ord}}\}$, thresholds for ordinal outcomes.

The outlined model implies the following likelihood based on observed data:

$$\begin{aligned} p(\mathbb{Y}_i | \mathcal{C}_i; \boldsymbol{\beta}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\tau}, \boldsymbol{\gamma}) &= \int \int p(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{\text{OB}}, \mathbb{Y}_i^{\text{*,OB}}, \mathbf{b}_i | \mathcal{C}_i; \boldsymbol{\beta}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\tau}, \boldsymbol{\gamma}) d\mathbf{b}_i d\mathbb{Y}_i^{\text{*,OB}} = \\ &= \int \int \underbrace{p(\mathbb{Y}_i^{\text{OB}} | \mathbb{Y}_i^{\text{*,OB}}; \boldsymbol{\gamma})}_{\text{thresholding (3)}} \cdot \underbrace{p(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{\text{*,OB}} | \mathcal{C}_i, \mathbf{b}_i; \boldsymbol{\beta}, \boldsymbol{\tau})}_{\text{MV LME (2),(4)}} \cdot \underbrace{p(\mathbf{b}_i | \boldsymbol{\mu}, \boldsymbol{\Sigma})}_{(5)} d\mathbf{b}_i d\mathbb{Y}_i^{\text{*,OB}}. \end{aligned} \quad (6)$$

The probability density functions which are integrated in (6) are of the form

$$\begin{aligned} p(\mathbb{Y}_i^{\text{OB}} | \mathbb{Y}_i^{\text{*,OB}}; \boldsymbol{\gamma}) &= \prod_{r \in \mathcal{R}^{\text{OB}}} \prod_{j=1}^{n_i} \left[\sum_{l=0}^{L^r-1} \mathbb{1}_{\{l\}}(y_{i,j}^r) \mathbb{1}_{(\gamma_l^r, \gamma_{l+1}^r]}(y_i^{\text{*,OB}}) \right], \\ p(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{\text{*,OB}} | \mathcal{C}_i, \mathbf{b}_i; \boldsymbol{\beta}, \boldsymbol{\tau}) &= \prod_{r \in \mathcal{R}^{\text{Num}}} \prod_{j=1}^{n_i} \varphi(y_{i,j}^r; \eta_{i,j}^r, \tau_r^{-1}) \cdot \prod_{r \in \mathcal{R}^{\text{OB}}} \prod_{j=1}^{n_i} \varphi(y_{i,j}^r; \eta_{i,j}^r, 1), \\ p(\mathbf{b}_i | \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \varphi(\mathbf{b}_i; \boldsymbol{\mu}, \boldsymbol{\Sigma}), \end{aligned} \quad (7)$$

where $\varphi(\cdot; \mathbf{m}, \mathbf{S})$ is probability density function of multivariate normal distribution with mean \mathbf{m} and variance matrix \mathbf{S} .

3 Model based clustering framework

Classification of subjects into one of K latent subgroups with apriori unknown structure will be based on the model-based clustering procedure developed above the model introduced in Section 2 in which all parameters of the underlying linear mixed models might be group specific. As it is usual in this context, let $U_i \in \{1, \dots, K\}$ denote an unobservable group-allocation indicator for subject i ($i = 1, \dots, n$). We assume that the model for i -th subject if it belongs to the k -th group (given $U_i = k, k = 1, \dots, K$) is described by the probability density function $p(\mathbb{Y}_i | \mathcal{C}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\gamma})$ of the form (6), where $\{\boldsymbol{\beta}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)}, \boldsymbol{\tau}^{(k)}\}$ is a set of (possibly) group-specific model parameters. That is, the assumed conditional probability distribution function of the i -th subject outcomes given the group allocation is

$$\begin{aligned} p(\mathbb{Y}_i | U_i = k, \mathcal{C}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\gamma}) &\stackrel{(6)}{=} \\ &= \int \int p(\mathbb{Y}_i^{\text{OB}} | \mathbb{Y}_i^{\text{*,OB}}; \boldsymbol{\gamma}) \cdot p(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{\text{*,OB}} | \mathcal{C}_i, \mathbf{b}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}) \cdot p(\mathbf{b}_i | \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)}) d\mathbf{b}_i d\mathbb{Y}_i^{\text{*,OB}}. \end{aligned} \quad (8)$$

Note that setting different LMM model parameters to be group-specific, we allow for different expressions of heterogeneity in the population. If, for example, we set parameters $\boldsymbol{\beta}$ to be group-specific we suppose that differences among the K latent groups can be described in terms of the effect of the fixed effects covariates \mathbb{X}_i . On the other hand, group-specific parameter $\boldsymbol{\Sigma}$ would lead to different associations among random effects that would subsequently change the marginal relationships among the outcomes. In general, not all of the LMM model parameters must be group-specific, nevertheless, for clarity, we suppress this in notation. In the following, symbols $\boldsymbol{\beta}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$ and $\boldsymbol{\tau}$ will represent sets of all corresponding parameters $\{\boldsymbol{\beta}^{(k)}, k = 1, \dots, K\}, \{\boldsymbol{\mu}^{(k)}, k = 1, \dots, K\}, \{\boldsymbol{\Sigma}^{(k)}, k = 1, \dots, K\}$ and $\{\boldsymbol{\tau}^{(k)}, k = 1, \dots, K\}$, respectively.

Let $w_k = P(U_i = k | \boldsymbol{w}) \in (0, 1), k = 1, \dots, K, \sum_{k=1}^K w_k = 1$, be (unknown) probabilities of pertinence to each of K groups, $\boldsymbol{w} := (w_1, \dots, w_K)$. Would we know all the model parameters $\boldsymbol{\theta} := \{\boldsymbol{w}, \boldsymbol{\beta}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\tau}, \boldsymbol{\gamma}\}$, Bayes rule provides an expression of conditional (given observed data) probabilities for subject i belonging to each of the groups:

$$u_{i,k}(\boldsymbol{\theta}) := P[U_i = k | \mathbb{Y}_i, \mathcal{C}_i; \boldsymbol{\theta}] = \frac{w_k p(\mathbb{Y}_i | U_i = k, \mathcal{C}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\gamma})}{\sum_{k'=1}^K w_{k'} p(\mathbb{Y}_i | U_i = k', \mathcal{C}_i; \boldsymbol{\beta}^{(k')}, \boldsymbol{\mu}^{(k')}, \boldsymbol{\Sigma}^{(k')}, \boldsymbol{\tau}^{(k')}, \boldsymbol{\gamma})}. \quad (9)$$

In a majority of the MBC methodologies, the authors consider maximum-likelihood estimation (MLE) of the unknown parameters. The clustering is then based on estimated subject specific group probabilities $\hat{u}_{i,k}^{\text{ML}} = u_{i,k}(\hat{\boldsymbol{\theta}}^{\text{ML}})$ where $\hat{\boldsymbol{\theta}}^{\text{ML}}$ denotes the MLE. In our situation, it would maximize the likelihood function:

$$L(\boldsymbol{\theta}) = \prod_{i=1}^n \left\{ \sum_{k=1}^K w_k p\left(\mathbb{Y}_i \mid U_i = k, \mathcal{C}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\gamma}\right) \right\},$$

This is traditionally solved by using the EM algorithm (Dempster et al, 1977) to face a problem of latent allocations leading to the mixture type likelihood. Nevertheless, two other levels of latent variables are present in our model leading to two additional levels of integration when evaluation the likelihood, see expression (6). They are the random effects \mathbf{b}_i and the latent numeric variables $\mathbb{Y}_i^{*,\text{OB}}$ associated with ordinal and binary outcomes \mathbb{Y}_i^{OB} as well. This makes the likelihood hardly tractable and we switch to the Bayesian framework and the related Markov chain Monte Carlo (MCMC) methodology which allows to fully exploit a hierarchical structure of our model. The clustering itself will then be based on the posterior distribution of the individual group probabilities (9).

4 Bayesian inference

For Bayesian inference, we exploit ideas of Bayesian Data Augmentation (BDA, Tanner and Wong, 1987) while considering all latent quantities, i.e., component allocations $\mathbf{U} := \{U_i, i = 1, \dots, n\}$, LMM random effect vectors $\mathbf{b} := \{\mathbf{b}_i, i = 1, \dots, n\}$ and latent variables $\mathbb{Y}^{*,\text{OB}} := \{\mathbb{Y}_i^{*,\text{OB}}, i = 1, \dots, n\}$ as additional model parameters included in the posterior distribution. With the model specified in Sections 2 and 3, the joint distribution of observed as well as latent data and model parameters for the Bayesian model is given by the following decomposition

$$\begin{aligned} p(\mathbb{Y}^{\text{N}}, \mathbb{Y}^{\text{OB}}, \mathbb{Y}^{*,\text{OB}}, \mathbf{U}, \mathbf{b}, \boldsymbol{\theta} \mid \mathcal{C}) &= \left[\prod_{i=1}^n p(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{\text{OB}}, \mathbb{Y}_i^{*,\text{OB}}, U_i, \mathbf{b}_i \mid \mathcal{C}_i; \boldsymbol{\theta}) \right] p(\boldsymbol{\theta}) \\ &= \left[\prod_{i=1}^n p(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{\text{OB}}, \mathbb{Y}_i^{*,\text{OB}} \mid \mathcal{C}_i, \mathbf{b}_i, U_i; \boldsymbol{\theta}) p(\mathbf{b}_i \mid U_i; \boldsymbol{\theta}) p(U_i \mid \boldsymbol{\theta}) \right] p(\boldsymbol{\theta}) \\ &= \left[\prod_{i=1}^n p(\mathbb{Y}_i^{\text{OB}} \mid \mathbb{Y}_i^{*,\text{OB}}; \boldsymbol{\gamma}) p(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{*,\text{OB}} \mid \mathcal{C}_i, \mathbf{b}_i; \boldsymbol{\beta}^{(U_i)}, \boldsymbol{\tau}^{(U_i)}) p(\mathbf{b}_i \mid \boldsymbol{\mu}^{(U_i)}, \boldsymbol{\Sigma}^{(U_i)}) w_{U_i} \right] p(\boldsymbol{\theta}), \end{aligned} \quad (10)$$

where factors in (10) follow from (7) and $p(\boldsymbol{\theta})$ is the prior distribution of the primary model parameters. By symbols \mathbb{Y}^{N} and \mathbb{Y}^{OB} we understand a collection of corresponding outcomes of the same type across all individuals, i.e. $\mathbb{Y}^{\text{N}} = \{\mathbb{Y}_i^{\text{N}}, i = 1, \dots, n\}$ and $\mathbb{Y}^{\text{OB}} = \{\mathbb{Y}_i^{\text{OB}}, i = 1, \dots, n\}$. Later, we also use \mathbb{Y} for all outcomes available, that is $\mathbb{Y} = \mathbb{Y}^{\text{N}} \cup \mathbb{Y}^{\text{OB}}$.

4.1 Prior distribution

We consider rather standard prior distributions of primary model parameters $\boldsymbol{\theta}$ being used in a context of hierarchical models being similar to that of ours. In particular, we assume that the prior distribution is decomposed as

$$p(\boldsymbol{\theta}) = p(\mathbf{w}) p(\boldsymbol{\gamma}) p(\boldsymbol{\beta} \mid \boldsymbol{\tau}) p(\boldsymbol{\tau}) p(\boldsymbol{\mu}) p(\boldsymbol{\Sigma})$$

with the following choices for the elements of the factorization.

A classically considered Dirichlet prior is assumed for the vector of group weights $\mathbf{w} = (w_1, \dots, w_K)$, i.e.,

$$p(\mathbf{w}) \propto \prod_{k=1}^K w_k^{\alpha_k - 1},$$

where $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_K)$ is a set of positive hyperparameters (all being equal to 1 in our applications in Sections 5 and 6).

Considering the thresholding parameters $\boldsymbol{\gamma}^r$, $r \in \mathcal{R}^{\text{Ord}}$ we first tackle the identifiability issue. Corresponding parametric space Ω^r is limited to set of all vectors of ordered values with fixed first threshold γ_1^r . An improper uniform distribution on Ω^r is assumed for each set of thresholds $\boldsymbol{\gamma}^r$, $r \in \mathcal{R}^{\text{Ord}}$. That is,

$$p(\boldsymbol{\gamma}) = \prod_{r \in \mathcal{R}^{\text{Ord}}} p(\boldsymbol{\gamma}^r) \propto \prod_{r \in \mathcal{R}^{\text{Ord}}} \mathbb{1}_{\Omega^r}(\boldsymbol{\gamma}^r).$$

All fixed effects parameters $\boldsymbol{\beta}^{r,(k)} = (\beta_1^{r,(k)}, \dots, \beta_{d_F^r}^{r,(k)})$, $r \in \mathcal{R}$, $k = 1, \dots, K$, are assumed to be apriori independent and following a conjugate normal distributions, i.e.,

$$p(\boldsymbol{\beta} | \boldsymbol{\tau}) = \prod_{k=1}^K \prod_{r \in \mathcal{R}^{\text{Num}}} \prod_{j=1}^{d_F^r} \varphi(\beta_j^{r,(k)}; \beta_{0,j}^r, (\tau_r^{(k)})^{-1} d_{j,j}^r) \cdot \prod_{k=1}^K \prod_{r \in \mathcal{R}^{\text{OB}}} \prod_{j=1}^{d_F^r} \varphi(\beta_j^{r,(k)}; \beta_{0,j}^r, d_{j,j}^r),$$

where $\beta_{0,j}^r$ and $d_{j,j}^r$ are fixed hyperparameters (being equal to zero and ten, respectively, in our applications). The precision parameters are given independent gamma priors, i.e.,

$$p(\boldsymbol{\tau}) = \prod_{k=1}^K \prod_{r \in \mathcal{R}} p(\tau_r^{(k)}),$$

where each $p(\tau_r^{(k)})$ corresponds to the gamma distribution $\Gamma(a_1, a_2)$.

Also for the random effect means, a set of independent, for simplicity only semi-conjugate normal priors is assumed, i.e.,

$$p(\boldsymbol{\mu}) \equiv p(\boldsymbol{\mu} | \boldsymbol{\tau}_R) = \prod_{k=1}^K \prod_{j=1}^{d_R} p(\mu_j^{(k)} | \tau_{R,j}^{(k)}) = \prod_{k=1}^K \prod_{j=1}^{d_R} \varphi(\mu_j^{(k)}; \mu_{0,j}^{(k)}, (\tau_{R,j}^{(k)})^{-1}),$$

where $\mu_{0,j}^{(k)}$ are fixed hyperparameters (being equal to zero in our applications). Finally, parameters $\boldsymbol{\tau}_R = \{\tau_R^{(k)}, k = 1, \dots, K\}$, $\boldsymbol{\tau}_R^{(k)} = (\tau_{R,1}^{(k)}, \dots, \tau_{R,d_R}^{(k)})$ are random hyperparameters being assigned independent gamma priors $\Gamma(a_3, a_4)$ in another level of hierarchy to allow for a weakly informative prior distribution.

Covariance matrices $\boldsymbol{\Sigma}^{(k)}$ of random effects \mathbf{b}_i are required to be completely general positive definite matrices, therefore, we suppose inverse covariance matrix $\boldsymbol{\Sigma}^{-(k)} := (\boldsymbol{\Sigma}^{(k)})^{-1}$ to follow Wishart distribution to preserve conjugacy. Again, to achieve a weakly informative prior we introduce a new hyperparameter, scale matrix $\mathbb{Q}^{(k)}$, while keeping the number of degrees of freedom $\nu_0 \geq d^R$ fixed. Inverse $\mathbb{Q}^{-(k)}$ of auxiliary scale matrix $\mathbb{Q}^{(k)}$ is also assumed to follow Wishart distribution, this time with fixed diagonal scale matrix $\mathbb{D}^{\mathbb{Q}}$ and number of degrees of freedom ν_1 . In our applications we use $\nu_0 = \nu_1 = d^R + 1$ and $\mathbb{D}^{\mathbb{Q}} = 100 \cdot \mathbb{I}_{d^R}$.

4.2 MCMC sampling scheme

Posterior distribution $p(\boldsymbol{\theta} | \mathbb{Y}, \mathcal{C})$ of model parameters $\boldsymbol{\theta}$ and its characteristics will be estimated using MCMC methodology Brooks et al (2011). We adopted well known Gibbs sampling scheme which samples a new value of each of the parameters from its full-conditioned distribution while always utilizing the last known values of other parameters. Due to our (semi)-conjugate choice of prior distributions the full-conditioned distributions are from well-known families, thus straightforward to be sampled from. Derivations of all full-conditioned distributions are postponed to Appendix section A. In this section we present only the set of required full-conditioned distributions for Gibbs sampling algorithm using the notation declared in Appendix:

1. $\mathbf{w} | \mathbf{U} \stackrel{(37)}{\sim} \text{Dir}_K(\mathbf{n}(\mathbf{U}) + \boldsymbol{\alpha})$,
2. $P(U_i = k | \mathbb{Y}_i^{\text{N}}, \mathcal{C}_i; \mathbb{Y}_i^{\text{OB}}, \mathbf{b}_i; \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{w}) \stackrel{(38)}{\equiv} \dots$,
3. $Y_{i,j}^{*,r} | Y_{i,j}^r = l, \boldsymbol{\gamma} \stackrel{(39)}{\sim} \text{TN}(\eta_{i,j}^{(U_i),r}, 1, \gamma_l^r, \gamma_{l+1}^r) \quad \text{for } r \in \mathcal{R}^{\text{OB}}, l = 0, \dots, L^r - 1$,
4. $\gamma_l^r | \mathbf{Y}^r; \mathbf{Y}^{*,r} \stackrel{(40)}{\sim} \text{Unif} \left[\max_{y \in \mathcal{Y}_{l-1}^r} y, \min_{y \in \mathcal{Y}_l^r} y \right] \quad \text{for } r \in \mathcal{R}^{\text{Ord}}, l = 2, \dots, L^r - 1$,
5. $\tau_r^{(k)} | \mathbf{Y}^r, \mathcal{C}^r; \mathbf{U}, \mathbf{b}^r; \boldsymbol{\beta}^{(k),r} \stackrel{(42)}{\sim} \Gamma(\tilde{a}_1^{(k),r}, \tilde{a}_2^{(k),r}) \quad \text{for } r \in \mathcal{R}^{\text{Num}}$,
6. $\boldsymbol{\beta}^{(k),r} | \mathbf{Y}^r, \mathcal{C}^r; \mathbf{U}, \mathbf{b}^r; \tau_r^{(k)} \stackrel{(44)}{\sim} N_{d_F^r}(\tilde{\boldsymbol{\beta}}^{(k),r}, (\tau_r^{(k)})^{-1} \left[(\mathbb{X}_{\mathcal{N}_k^r}^r(\mathbf{U}))^\top \mathbb{X}_{\mathcal{N}_k^r}^r(\mathbf{U}) + (\mathbb{D}^r)^{-1} \right]^{-1}) \quad \text{for } r \in \mathcal{R}$,
7. $\tau_{R,j}^{(k)} | \mu_j^{(k)} \stackrel{(46)}{\sim} \Gamma\left(a_3 + \frac{1}{2}, a_4 + \frac{1}{2} (\mu_j^{(k)} - \mu_{0,j}^{(k)})^2\right)$,
8. $\boldsymbol{\mu}^{(k)} | \mathbf{U}, \mathbf{b}; \boldsymbol{\Sigma}^{(k)}; \boldsymbol{\tau}_R^{(k)} \stackrel{(48)}{\sim} N_{d^R}(\tilde{\boldsymbol{\mu}}^{(k)}, [n^k(\mathbf{U})\boldsymbol{\Sigma}^{-(k)} + \text{diag}(\boldsymbol{\tau}_R^{(k)})]^{-1})$,

9. $\mathbb{Q}^{-(k)} \mid \boldsymbol{\Sigma}^{(k)} \stackrel{(50)}{\sim} W_{dR} \left(\left[\boldsymbol{\Sigma}^{-(k)} + (\mathbb{D}^{\mathbb{Q}})^{-1} \right]^{-1}, \mathbf{v}_0 + \mathbf{v}_1 \right),$
10. $\boldsymbol{\Sigma}^{-(k)} \mid \mathbf{U}, \mathbf{b}; \boldsymbol{\mu}^{(k)}, \mathbb{Q}^{(k)} \stackrel{(52)}{\sim} W_{dR} \left(\tilde{\mathbb{Q}}^{(k)}, n^k(\mathbf{U}) + \mathbf{v}_0 \right),$
11. $\mathbf{b}_i \mid \mathbb{Y}_i^N, \mathcal{C}_i; \mathbb{Y}_i^{*,OB}, U_i; \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\mu}, \boldsymbol{\Sigma} \stackrel{(54)}{\sim} N_{dR} \left(\tilde{\mathbf{b}}_i, \left[\tilde{\mathbb{Z}}_i^\top \tilde{\mathbb{Z}}_i + \boldsymbol{\Sigma}^{-(k)} \right]^{-1} \right),$

A sample from converged chain constructed in this way then can be used for estimation of marginal posterior distributions of individual parameters.

4.3 Classification probabilities

Posterior distribution is invariant towards permutation of cluster labels. Hence, before calculation of classification probabilities we consider problem of *label switching* by method of Stephens (2000). This post-sampling procedure that considers all $K!$ permutations of labels for each iteration ensures that the latent clusters $1, \dots, K$ have fixed meaning during the whole sampling procedure. However, in our applications no need of permutation has occurred.

Primarily, we perform classification using posterior means $\hat{U}_{i,k} = \int_{\boldsymbol{\theta}} u_{i,k}(\boldsymbol{\theta}) \cdot p(\boldsymbol{\theta} \mid \mathbb{Y}, \mathcal{C}) d\boldsymbol{\theta}$ of allocation probabilities $u_{i,k}(\boldsymbol{\theta})$ defined in (9). With MCMC based inference we need to calculate $u_{i,k}(\boldsymbol{\theta})$ at each iteration of sampled $\boldsymbol{\theta}$, for which probability density functions of i -th outcomes given cluster need to be evaluated (8). Unfortunately, this involves non-trivial integration of auxiliary latent variables - random effects \mathbf{b}_i and latent numeric outcomes $\mathbb{Y}_i^{*,OB}$. Therefore, we devote the rest of this section to description of a method chosen for calculating desired integrals

$$\begin{aligned}
p\left(\mathbb{Y}_i \mid U_i = k, \mathcal{C}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\gamma}\right) &= \\
&= \int p\left(\mathbb{Y}_i^{OB} \mid \mathbb{Y}_i^{*,OB}; \boldsymbol{\gamma}\right) \cdot \underbrace{\left[\int p\left(\mathbb{Y}_i^N, \mathbb{Y}_i^{*,OB} \mid \mathcal{C}_i, \mathbf{b}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}\right) \cdot p\left(\mathbf{b}_i \mid \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)}\right) d\mathbf{b}_i \right]}_{p\left(\mathbb{Y}_i^N, \mathbb{Y}_i^{*,OB} \mid \mathcal{C}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)}\right)} d\mathbb{Y}_i^{*,OB}. \quad (11)
\end{aligned}$$

4.3.1 Integration with respect to random effects \mathbf{b}_i

Let us first perform the integration of random effects \mathbf{b}_i from (11) to obtain distribution of numeric variables unconditioned by random effects. We will avoid integration by realization that under normality assumption of both numeric outcomes and random effects the unconditioned distribution of outcomes is also normal.

Vector of all numeric and latent numeric outcomes \mathbf{Y}_i (\mathbb{Y}_i^N combined with $\mathbb{Y}_i^{*,OB}$) of length $d = n_i \cdot |\mathcal{R}|$ given a vector of all random effects \mathbf{b}_i follows by our LME assumption multivariate normal distribution

$$\mathbf{Y}_i \mid \mathcal{C}_i; \mathbf{b}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)} \sim N_d \left(\mathbb{X}_i \boldsymbol{\beta}^{(k)} + \mathbb{Z}_i \mathbf{b}_i, \mathbb{T} \right),$$

where \mathbb{X}_i and \mathbb{Z}_i are block diagonal matrices composed of model matrices of fixed effects \mathbb{X}_i^f and of random effects \mathbb{Z}_i^r , respectively. The variance matrix \mathbb{T} is diagonal due to independence assumption and contains corresponding variance, that is τ_r^{-1} for $r \in \mathcal{R}^{\text{Num}}$ and 1 otherwise.

Using the normality of random effects, i.e. $\mathbf{b}_i \mid \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)} \sim N_{dR} \left(\boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)} \right)$, and notoriously known formulas we obtain unconditioned mean and variance matrix:

$$\begin{aligned}
\mathbb{E} \left[\mathbf{Y}_i \mid \mathcal{C}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)} \right] &= \mathbb{E}(\mathbb{E}[\mathbf{Y}_i \mid \mathbf{b}_i]) = \mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \boldsymbol{\mu}, \\
\text{var} \left[\mathbf{Y}_i \mid \mathcal{C}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)} \right] &= \mathbb{E}(\text{var}[\mathbf{Y}_i \mid \mathbf{b}_i]) + \text{var}(\mathbb{E}[\mathbf{Y}_i \mid \mathbf{b}_i]) = \mathbb{T} + \mathbb{Z}^\top \boldsymbol{\Sigma} \mathbb{Z} =: \mathbb{V},
\end{aligned}$$

which results in

$$\mathbf{Y}_i \mid \mathcal{C}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)} \sim N_d(\mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \boldsymbol{\mu}, \mathbb{V}). \quad (12)$$

This distribution has general covariance structure which reflects the general structure of $\boldsymbol{\Sigma}$ and hence captures dependencies among outcomes.

4.3.2 Integration with respect to latent numeric outcomes $\mathbb{Y}_i^{*,\text{OB}}$

It remains to perform the following integration:

$$\int p\left(\mathbb{Y}_i^{\text{OB}} \mid \mathbb{Y}_i^{*,\text{OB}}; \boldsymbol{\gamma}\right) \cdot p\left(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{*,\text{OB}} \mid \mathcal{C}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)}\right) d\mathbb{Y}_i^{*,\text{OB}},$$

which is in fact integration of multivariate normal density within the bounds given by thresholds $\boldsymbol{\gamma}$ and observed ordinal and binary outcomes. First, we separate marginal distribution of numeric outcomes \mathbb{Y}_i^{N} since it can avoid integration, while the conditional normal distribution of latent numeric outcomes $\mathbb{Y}_i^{*,\text{OB}}$ given \mathbb{Y}_i^{N} still awaits the integration:

$$\underbrace{p\left(\mathbb{Y}_i^{\text{N}} \mid \mathcal{C}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)}\right)}_{\text{pdf of MVN}} \cdot \underbrace{\int p\left(\mathbb{Y}_i^{\text{OB}} \mid \mathbb{Y}_i^{*,\text{OB}}; \boldsymbol{\gamma}\right)}_{\text{thresholding (3)}} \cdot \underbrace{\varphi\left(\mathbb{Y}_i^{*,\text{OB}}; \boldsymbol{\eta}_{\text{OB}}^{(k)}, \mathbb{V}_{\text{OB}}^{(k)}\right)}_{\text{pdf of } \mathbb{Y}_i^{*,\text{OB}} \mid \mathbb{Y}_i^{\text{N}}} d\mathbb{Y}_i^{*,\text{OB}},$$

where $\boldsymbol{\eta}_{\text{OB}}^{(k)}$ and $\mathbb{V}_{\text{OB}}^{(k)}$ are the conditional mean and variance matrix of $\mathbb{Y}_i^{*,\text{OB}} \mid \mathbb{Y}_i^{\text{N}}, \mathcal{C}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)}$.

It remains to integrate the product of two functions, first one of which just declares lower and upper integration bounds and the second one is the probability density function of multivariate normal distribution with mean $\boldsymbol{\eta}_{\text{OB}}^{(k)}$ and variance matrix $\mathbb{V}_{\text{OB}}^{(k)}$. For each individual categorical outcome $r \in \mathcal{R}^{\text{OB}}$ and observation $j \in \{1, \dots, n_i\}$ the value $y_{i,j}^r = l$ determines an interval given by the corresponding pair of $\boldsymbol{\gamma}$ parameters, see (3):

$$Y_{i,j}^r = l \implies Y_{i,j}^{*,r} \in (\boldsymbol{\gamma}_l^r, \boldsymbol{\gamma}_{l+1}^r] =: (e_{ij}^r, f_{ij}^r].$$

If we denote the resulting Cartesian product of these intervals as $\square(\mathbb{Y}_i^{\text{OB}}) = (\mathbf{e}_i, \mathbf{f}_i] \subset \mathbb{R}^{d^{\text{OB}}}$ then the remaining integral can be written in the form

$$I_k(\square) = \int_{\square(\mathbb{Y}_i^{\text{OB}})} p\left(\mathbb{Y}_i^{*,\text{OB}} \mid \mathbb{Y}_i^{\text{N}}, \mathcal{C}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)}\right) d\mathbb{Y}_i^{*,\text{OB}} = \int_{\mathbf{e}_i}^{\mathbf{f}_i} \varphi\left(\mathbf{y}; \boldsymbol{\eta}_{\text{OB}}^{(k)}, \mathbb{V}_{\text{OB}}^{(k)}\right) d\mathbf{y}. \quad (13)$$

Finally, after the integrals I_k for all $k = 1, \dots, K$ are computed, the classification probabilities can be calculated proportionally:

$$u_{i,k}(\boldsymbol{\theta}) = \frac{w_k \cdot p\left(\mathbb{Y}_i^{\text{N}} \mid \mathcal{C}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)}\right) \cdot I_k(\square)}{\sum_{k'=1}^K w_{k'} \cdot p\left(\mathbb{Y}_i^{\text{N}} \mid \mathcal{C}_i; \boldsymbol{\beta}^{(k')}, \boldsymbol{\tau}^{(k')}, \boldsymbol{\mu}^{(k')}, \boldsymbol{\Sigma}^{(k')}\right) \cdot I_{k'}(\square)}. \quad (14)$$

In order to compute integrals (13) needed in (14) we adopted an effective algorithm presented by Genz (1992) which is also based on MCMC sampling. Since the approximation of such an integral is needed K -times for each generated state of our Gibbs sampling, the procedure is considerably time-consuming. The implemented function `pvmnorm` in  package `mvtnorm` (Genz et al, 2019) is being used in our applications.

4.4 Classification rules

Once we have estimated the posterior means of classification probabilities $\widehat{U}_{i,k}$ we can face the problem of pairing a subject i with the most suitable cluster. Naturally, we may choose cluster k such that the corresponding estimated $\widehat{U}_{i,k}$ is the largest among all $k = 1, \dots, K$ values. However, that may not be the most fitting choice in cases when two of the clusters have both comparable and high probability.

In order to prevent misclassification, we allow subjects to remain unclassified when the decision is not clear. One way to accomplish that would be to classify into the cluster with the highest probability only if it clearly overcomes the second largest probability. That is, when the difference between the two largest probabilities is higher than chosen threshold. However, the choice of the value of this threshold for different values of K would be another problem to be dealt with. Therefore, in our applications we make use of 95% Highest Posterior Density (HPD) interval estimates. We classify a subject i into the class k with the highest classification probability $\widehat{U}_{i,k}$ if its lower 95% HPD bound is still higher than any other upper 95% HPD bound of the remaining probabilities. Otherwise, the subject i remained unclassified. That should fill clusters with their most typical representatives and keep indecisive subjects aside. Unclassified subjects can then be additionally analysed to determine the pair (or potentially larger group) of clusters they are associated with the most.

4.5 Number of groups

Throughout the paper we treated the number of latent classes K fixed as to be selected by statistician in advance. In most circumstances, however, there is no prior knowledge of the suitable value of K to be used. Usual practice in this situation is to try several values of K and choose the one that optimizes one of the known criteria. We follow the steps of [Aitkin et al \(2009\)](#) and use a deviance criterion that is solely based on the goodness of fit measured by the log-likelihood function.

Deviance is a parametric function of $\boldsymbol{\theta}$ generally defined as

$$D(\boldsymbol{\theta}; \mathbb{Y}, \mathcal{C}) = -2 \log p(\mathbb{Y} | \mathcal{C}; \boldsymbol{\theta}) = -2 \sum_{i=1}^n \log p(\mathbb{Y}_i | \mathcal{C}_i; \boldsymbol{\theta}). \quad (15)$$

Contribution of individual i to the deviance in our model can be expressed as

$$-2 \log p(\mathbb{Y}_i | \mathcal{C}_i; \boldsymbol{\theta}) = -2 \log \left[\sum_{k=1}^K \int \int p \left(\mathbb{Y}_i, \mathbb{Y}_i^{*, \text{OB}}, \mathbf{b}_i, U_i = k \mid \mathcal{C}_i; \boldsymbol{\theta} \right) d\mathbf{b}_i d\mathbb{Y}_i^{*, \text{OB}} \right] \quad (16)$$

which includes the integration (11) for calculating classification probabilities which has been described in section 4.7. Using the same notation we can write

$$D(\boldsymbol{\theta}; \mathbb{Y}, \mathcal{C}) = -2 \sum_{i=1}^n \log \left[\sum_{k=1}^K w_k \cdot p \left(\mathbb{Y}_i^{\text{N}} \mid \mathcal{C}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)} \right) \cdot I_k \left(\square \left(\mathbb{Y}_i^{\text{OB}} \right) \right) \right] \quad (17)$$

where the denominator of (14) is inserted into the logarithm. Therefore, calculation of deviance for one set of parameters $\boldsymbol{\theta}$ requires calculation of classification probabilities for every individual. Only then the deviance parametric function can be fully evaluated. Hence, exploration of posterior distribution of deviance is heavily time-consuming. This is why we applied thinning in our applications to speed up calculations of all of the probabilities.

4.6 MCMC sampling scheme

Joint density of observed outcomes \mathbb{Y}_i of i th subject and its latent variables U_i , $\mathbb{Y}_i^{*,\text{OB}}$ and \mathbf{b}_i decomposes into product of several conditional probability density functions:

$$p\left(\mathbb{Y}_i, U_i, \mathbb{Y}_i^{*,\text{OB}}, \mathbf{b}_i \mid \mathcal{C}_i; \boldsymbol{\theta}\right) = p\left(\mathbb{Y}_i^{\text{OB}} \mid \mathbb{Y}_i^{*,\text{OB}}; \boldsymbol{\gamma}\right) \cdot p\left(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{*,\text{OB}} \mid \mathcal{C}_i, \mathbf{b}_i; \boldsymbol{\beta}^{(U_i)}, \boldsymbol{\tau}^{(U_i)}\right) \cdot p\left(\mathbf{b}_i \mid \boldsymbol{\mu}^{(U_i)}, \boldsymbol{\Sigma}^{-1(U_i)}\right) \cdot p(U_i \mid \mathbf{w}), \quad (18)$$

which is easier to work with than $p(\mathbb{Y}_i \mid \mathcal{C}_i; \boldsymbol{\theta})$ obtainable by combining (??) and (8). We adopt Bayesian framework in which we consider parameters $\boldsymbol{\theta}$ to be of random nature by assigning a prior distribution $p(\boldsymbol{\theta})$ of parameter $\boldsymbol{\theta}$. Then complete probability density function of all measured or latent data (across all subjects $i = 1, \dots, n$) and parameters is of the form:

$$p\left(\mathbb{Y}, \mathbf{U}, \mathbb{Y}^{*,\text{OB}}, \mathbf{b}, \boldsymbol{\theta} \mid \mathcal{C}\right) = \left[\prod_{i=1}^n p\left(\mathbb{Y}_i, U_i, \mathbb{Y}_i^{*,\text{OB}}, \mathbf{b}_i \mid \mathcal{C}_i; \boldsymbol{\theta}\right) \right] \cdot p(\boldsymbol{\theta}),$$

to which is posterior probability density function $p(\mathbf{U}, \mathbb{Y}^{*,\text{OB}}, \mathbf{b}, \boldsymbol{\theta} \mid \mathbb{Y}, \mathcal{C})$ proportional by the Bayes theorem. Ideally, we would be interested in marginal posterior distribution of parameter $\boldsymbol{\theta}$ given by

$$p(\boldsymbol{\theta} \mid \mathbb{Y}, \mathcal{C}) = \int \int \int p\left(\mathbf{U}, \mathbb{Y}^{*,\text{OB}}, \mathbf{b}, \boldsymbol{\theta} \mid \mathbb{Y}, \mathcal{C}\right) d\mathbf{b} d\mathbb{Y}^{*,\text{OB}} d\mathbf{U}.$$

To avoid computation of needed multiplicative constant and integration with respect to all latent variables we exploit MCMC methods (Robert, 2001) based on sampling Markov chain of states corresponding to parameter $\boldsymbol{\theta}$ and latent variables. As the chain reaches its stationary distribution (equal to the posterior one) generated states are considered to be sampled representatives of that distribution. Marginal posterior distribution of $\boldsymbol{\theta}$ can then be described by M generated values $\boldsymbol{\theta}_{[B+1]}, \dots, \boldsymbol{\theta}_{[B+M]}$ after the burn-in period B .

For our model we decided to choose basic Gibbs sampling method (Geman and Geman, 1984) for constructing Markov chain of required properties. Such an algorithm requires effortless sampling from so-called *full-conditioned* distributions. We can obtain a well-known distribution when prior distributions of all components of $\boldsymbol{\theta}$ are suitably set to achieve conjugacy.

4.6.1 Prior distributions

decomposed into several independently distributed blocks: \mathbf{w} , $\boldsymbol{\gamma}$, $(\boldsymbol{\beta}, \boldsymbol{\tau})$, $(\boldsymbol{\mu}, \boldsymbol{\tau}_R)$, $(\boldsymbol{\Sigma}^{-1}, \mathbb{Q}^{-1})$, which results in the following form of the probability density function:

$$p\left(\boldsymbol{\theta}, \mathbb{Q}^{-1}, \boldsymbol{\tau}_R \mid \boldsymbol{\zeta}\right) = p(\mathbf{w} \mid \boldsymbol{\zeta}) \cdot p(\boldsymbol{\gamma} \mid \boldsymbol{\zeta}) \cdot p(\boldsymbol{\beta} \mid \boldsymbol{\tau}; \boldsymbol{\zeta}) \cdot p(\boldsymbol{\tau} \mid \boldsymbol{\zeta}) \cdot p(\boldsymbol{\mu} \mid \boldsymbol{\tau}_R; \boldsymbol{\zeta}) \cdot p(\boldsymbol{\tau}_R \mid \boldsymbol{\zeta}) \cdot p(\boldsymbol{\Sigma}^{-1} \mid \mathbb{Q}^{-1}; \boldsymbol{\zeta}) \cdot p(\mathbb{Q}^{-1} \mid \boldsymbol{\zeta}),$$

where $\boldsymbol{\tau}^R$, \mathbb{Q}^{-1} are another random parameters (or set of parameters in case of group-specificity) making prior distributions of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}^{-1}$ more flexible and $\boldsymbol{\zeta}$ stands for set of other hyper-parameters that need to be fixed.

Since \mathbf{w} is a vector of probabilities that add up to 1, ideal choice for prior distribution is Dirichlet distribution $\text{Dir}_K(\boldsymbol{\alpha})$ with pdf equal to

$$p(\mathbf{w} \mid \boldsymbol{\alpha}) = \frac{\Gamma\left(\sum_{k=1}^K \alpha_k\right)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K w_k^{\alpha_k - 1}, \quad (19)$$

where $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_K)$ is a set of positive parameters (set to 1 by default).

For each ordinal variable $r^O \in \mathcal{R}^O$ we suppose that thresholding constants $\boldsymbol{\gamma}^{r^O}$ follow (independently) an improper uniform distribution on $\mathbb{R}_{\boldsymbol{\gamma}_1^O <}^{L^O - 2}$, i.e.

$$p\left(\boldsymbol{\gamma} \mid \boldsymbol{\gamma}_1^O, r^O \in \mathcal{R}^{\text{Ord}}\right) \propto \prod_{r^O \in \mathcal{R}^{\text{Ord}}} \mathbb{1}_{\mathbb{R}_{\boldsymbol{\gamma}_1^O <}^{L^O - 2}}\left(\boldsymbol{\gamma}^{r^O}\right). \quad (20)$$

For each outcome $r \in \mathcal{R}$ we suppose that fixed effects

$$\boldsymbol{\beta}^{(k),r} \mid \boldsymbol{\tau}_r^{(k)}; \boldsymbol{\beta}_0^r, \mathbb{D}^r \sim N_{d_r^F} \left(\boldsymbol{\beta}_0^r, \left(\boldsymbol{\tau}_r^{(k)} \right)^{-1} \mathbb{D}^r \right) \quad \text{for all } k = 1, \dots, K, \quad (21)$$

where $\boldsymbol{\beta}_0^r \in \mathbb{R}^{d_r^F}$ is the prior mean and \mathbb{D}^r is diagonal matrix with diagonal elements $d_{j,j}^r > 0$. We remind that for ordinal or binary variable we impose $\tau_{r,\text{OB}} = 1$ for $r^{\text{OB}} \in \mathcal{R}^{\text{OB}}$, which eliminates $\tau_{r,\text{OB}}$ from (21). Due to supposed prior independence of $\boldsymbol{\beta}^{(k),r}$ we get the following pdf:

$$p(\boldsymbol{\beta} \mid \boldsymbol{\tau}; \boldsymbol{\beta}_0, \mathbb{D}) = \prod_{k=1}^K \prod_{r \in \mathcal{R}} \prod_{j=1}^{d_r^F} \sqrt{\frac{\tau_r^{(k)}}{2\pi d_{j,j}^r}} \exp \left\{ -\frac{\tau_r^{(k)}}{2d_{j,j}^r} \left(\beta_j^{(k),r} - \beta_{0,j}^r \right)^2 \right\}. \quad (22)$$

For numerical outcomes $r^{\text{N}} \in \mathcal{R}^{\text{Num}}$ we set precision parameters $\tau_{r^{\text{N}}}^{(k)}$ to independently follow $\Gamma(a_1, a_2)$ with $a_1, a_2 > 0$ fixed hyper-parameters, which leads to

$$p(\boldsymbol{\tau} \mid a_1, a_2) = \prod_{k=1}^K \prod_{r^{\text{N}} \in \mathcal{R}^{\text{Num}}} \frac{a_2^{a_1}}{\Gamma(a_1)} \left(\tau_{r^{\text{N}}}^{(k)} \right)^{a_1-1} \exp \left\{ -a_2 \tau_{r^{\text{N}}}^{(k)} \right\}. \quad (23)$$

We set $\beta_{0,j}^r = 0$, $d_{j,j}^r = 10$ and $a_1 = a_2 = 1$.

For set of random effects $\boldsymbol{\mu}$ we suppose $\boldsymbol{\mu}^{(k)} \mid \boldsymbol{\tau}_R^{(k)}; \boldsymbol{\mu}_0 \sim N_{d^R} \left(\boldsymbol{\mu}_0, \text{diag} \left(\boldsymbol{\tau}_R^{(k)} \right)^{-1} \right)$ independently for all $k = 1, \dots, K$, where $\boldsymbol{\mu}_0 \in \mathbb{R}^{d^R}$ is the prior mean of $\boldsymbol{\mu}^{(k)}$ and its prior covariance matrix is diagonal with positive elements $\boldsymbol{\tau}_R^{(k)} = \left(\tau_{R,1}^{(k)}, \dots, \tau_{R,d^R}^{(k)} \right)^\top$. These diagonal elements (viewed as another random parameters of our hierarchical model) might not be considered to be group-specific, nevertheless, we from now on work under this assumption for the sake of generality. Similarly as for elements of $\boldsymbol{\tau}$ we suppose elements of $\boldsymbol{\tau}_R$ follow independently $\Gamma(a_3, a_4)$ with fixed hyper-parameters $a_3, a_4 > 0$, which results in probability density functions:

$$p(\boldsymbol{\mu} \mid \boldsymbol{\tau}_R; \boldsymbol{\mu}_0) = \prod_{k=1}^K \prod_{j=1}^{d^R} \sqrt{\frac{\tau_{R,j}^{(k)}}{2\pi}} \exp \left\{ -\frac{\tau_{R,j}^{(k)}}{2} \left(\mu_j^{(k)} - \mu_{0,j} \right)^2 \right\}, \quad (24)$$

$$p(\boldsymbol{\tau}_R \mid a_3, a_4) = \prod_{k=1}^K \prod_{j=1}^{d^R} \frac{a_4^{a_3}}{\Gamma(a_3)} \left(\tau_{R,j}^{(k)} \right)^{a_3-1} \exp \left\{ -a_4 \tau_{R,j}^{(k)} \right\}.$$

We set $\mu_{0,j}^{(k)} = 0$ and $a_3 = a_4 = 1$.

Covariance matrices $\boldsymbol{\Sigma}^{(k)}$ of random effects \boldsymbol{b}_i are required to be completely general positive definite matrices, therefore, we suppose inverse covariance matrix $\boldsymbol{\Sigma}^{-(k)}$ to follow Wishart distribution to preserve conjugacy. More specifically, $\boldsymbol{\Sigma}^{-(k)} \mid \mathbb{Q}^{-(k)}, \nu_0 \sim W_{d^R} \left(\mathbb{Q}^{-(k)}, \nu_0 \right)$, where $\mathbb{Q}^{(k)}$ is (for generality purposes) group-specific scale matrices and $\nu_0 \geq d^R$ is the number of degrees of freedom. Auxiliary positive definite random matrix $\mathbb{Q}^{-(k)}$ is distributed in very similar way: $\mathbb{Q}^{-(k)} \mid \mathbb{D}^{\mathbb{Q}}, \nu_1 \sim W_{d^R} \left(\mathbb{D}^{\mathbb{Q}}, \nu_1 \right)$, where $\mathbb{D}^{\mathbb{Q}}$ is diagonal matrix with fixed diagonal elements $d_{j,j}^{\mathbb{Q}} > 0$ and ν_1 is again the number of degrees of freedom. Probability density functions of $\boldsymbol{\Sigma}^{-1} = \left\{ \boldsymbol{\Sigma}^{-(k)}, k = 1, \dots, K \right\}$ and $\mathbb{Q}^{-1} = \left\{ \mathbb{Q}^{-(k)}, k = 1, \dots, K \right\}$ then can be expressed as

$$p(\boldsymbol{\Sigma}^{-1} \mid \mathbb{Q}^{-1}; \nu_0) \propto \prod_{k=1}^K \left| \mathbb{Q}^{-(k)} \right|^{\frac{\nu_0}{2}} \left| \boldsymbol{\Sigma}^{-(k)} \right|^{\frac{\nu_0 - d^R - 1}{2}} \exp \left\{ -\text{Tr} \left[\mathbb{Q}^{-(k)} \boldsymbol{\Sigma}^{-(k)} \right] \right\}, \quad (25)$$

$$p(\mathbb{Q}^{-1} \mid \mathbb{D}^{\mathbb{Q}}, \nu_1) \propto \prod_{k=1}^K \left| \mathbb{Q}^{-(k)} \right|^{\frac{\nu_1 - d^R - 1}{2}} \exp \left\{ -\text{Tr} \left[\left(\mathbb{D}^{\mathbb{Q}} \right)^{-1} \mathbb{Q}^{-(k)} \right] \right\}.$$

We set $\nu_0 = \nu_1 = d^R + 1$ and $\mathbb{D}^{\mathbb{Q}} = 100 \cdot \mathbb{I}_{d^R}$.

The set of all hyper-parameters then consists of

$$\boldsymbol{\zeta} = \left\{ \boldsymbol{\alpha}, \gamma_1^{\text{OB}}, \boldsymbol{\beta}_0^r, \mathbb{D}^r, a_1, a_2, \boldsymbol{\mu}_0, a_3, a_4, \nu_0, \nu_1, \mathbb{D}^{\mathbb{Q}}, r^{\text{OB}} \in \mathcal{R}^{\text{OB}}, r \in \mathcal{R} \right\}.$$

4.6.2 Gibbs algorithm

For construction of Gibbs algorithm it is needed to derive full-conditioned distributions of all random parameters that enter our model $\Psi = \{\mathbf{U}, \mathbb{Y}^{*,\text{OB}}, \mathbf{b}, \mathbf{w}, \boldsymbol{\gamma}, \boldsymbol{\beta}, \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}, \boldsymbol{\tau}, \boldsymbol{\tau}_R, \mathbb{Q}^{-1}\}$. By full-conditioned distribution of parameter $\boldsymbol{\psi} \in \Psi$ we understand distribution $\boldsymbol{\psi} | \mathbb{Y}, \mathcal{C}; \Psi_{-\boldsymbol{\psi}}; \boldsymbol{\zeta}$, where $\Psi_{-\boldsymbol{\psi}} = \Psi \setminus \{\boldsymbol{\psi}\}$, with probability density function

$$p(\boldsymbol{\psi} | \mathbb{Y}, \mathcal{C}; \Psi_{-\boldsymbol{\psi}}; \boldsymbol{\zeta}) \propto \left[\prod_{i=1}^n p\left(\mathbb{Y}_i, U_i, \mathbb{Y}_i^{*,\text{OB}}, \mathbf{b}_i | \mathcal{C}_i; \boldsymbol{\theta}\right) \right] \cdot p(\boldsymbol{\theta}, \mathbb{Q}^{-1}, \boldsymbol{\tau}_R | \boldsymbol{\zeta}) \quad (26)$$

viewed as a function of $\boldsymbol{\psi}$.

Derivations of all full-conditioned distributions of $\boldsymbol{\psi} \in \Psi$ are postponed to [Appendix section A](#). In this section we present only the resulting Gibbs sampling algorithm using the notation declared in [Appendix](#).

Algorithm By index $\bullet_{[m]}$ we understand value of parameter \bullet generated in m -th step.

1. Take initial values $\boldsymbol{\Psi}_{[0]}$ and set $m = 1$.
2. Generate m -th set of parameters $\boldsymbol{\Psi}_{[m]}$:
 - (a) $\mathbf{w} | \mathbf{U}; \boldsymbol{\alpha} \stackrel{(37)}{\sim} \text{Dir}_K(\mathbf{n}(\mathbf{U}) + \boldsymbol{\alpha})$,
 - (b) $P\left(U_i = k | \mathbb{Y}_i^N, \mathcal{C}_i; \mathbb{Y}_i^{*,\text{OB}}, \mathbf{b}_i; \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}, \mathbf{w}\right) \stackrel{(38)}{=} \dots$,
 - (c) $Y_{i,j}^{*,r\text{OB}} | Y_{i,j}^{r\text{OB}} = l, \boldsymbol{\gamma} \stackrel{(39)}{\sim} \text{TN}\left(\eta_{i,j}^{(U_i),r\text{OB}}, 1, \gamma_l^{r\text{OB}}, \gamma_{l+1}^{r\text{OB}}\right)$,
 - (d) $\gamma_l^{r\text{O}} | \mathbf{Y}^{r\text{O}}; \mathbf{Y}^{*,r\text{O}} \stackrel{(40)}{\sim} \text{Unif}\left[\max_{y \in \mathcal{Z}_{l-1}^{r\text{O}}} y, \min_{y \in \mathcal{Z}_l^{r\text{O}}} y\right]$, $l = 1, \dots, L^{r\text{O}}$,
 - (e) $\tau_{r,N}^{(k)} | \mathbf{Y}^{rN}, \mathcal{C}^{rN}; \mathbf{U}, \mathbf{b}^{rN}; \boldsymbol{\beta}^{(k),rN}; \boldsymbol{\beta}_0^{rN}, \mathbb{D}^{rN}, a_1, a_2 \stackrel{(42)}{\sim} \Gamma\left(\tilde{a}_1^{(k),rN}, \tilde{a}_2^{(k),rN}\right)$,
 - (f) $\boldsymbol{\beta}^{(k),r} | \mathbf{Y}^r, \mathcal{C}^r; \mathbf{U}, \mathbf{b}^r; \tau_r^{(k)}; \boldsymbol{\beta}_0^r, \mathbb{D}^r \stackrel{(44)}{\sim} N_{d_F^r}\left(\tilde{\boldsymbol{\beta}}^{(k),r}, \frac{\left[\left(\mathbb{X}_{\mathcal{N}_k^r}^r\right)^\top \mathbb{X}_{\mathcal{N}_k^r}^r(\mathbf{U}) + (\mathbb{D}^r)^{-1}\right]^{-1}}{\tau_r^{(k)}}\right)$,
 - (g) $\tau_{R,j}^{(k)} | \mu_j^{(k)}; \mu_{0,j}^{(k)}, a_3, a_4 \stackrel{(46)}{\sim} \Gamma\left(a_3 + \frac{1}{2}, a_4 + \frac{1}{2} \left(\mu_j^{(k)} - \mu_{0,j}^{(k)}\right)^2\right)$,
 - (h) $\boldsymbol{\mu}^{(k)} | \mathbf{U}, \mathbf{b}; \boldsymbol{\Sigma}^{-(k)}; \boldsymbol{\tau}_R^{(k)}, \boldsymbol{\mu}_0^{(k)} \stackrel{(48)}{\sim} N_{d_R}\left(\tilde{\boldsymbol{\mu}}^{(k)}, \left[n^k(\mathbf{U})\boldsymbol{\Sigma}^{-(k)} + \text{diag}\left(\boldsymbol{\tau}_R^{(k)}\right)\right]^{-1}\right)$,
 - (i) $\mathbb{Q}^{-(k)} | \boldsymbol{\Sigma}^{-(k)}; \mathbf{v}_0, \mathbf{v}_1, \mathbb{D}^{\mathbb{Q}} \stackrel{(50)}{\sim} W_{d_R}\left(\left[\boldsymbol{\Sigma}^{-(k)} + (\mathbb{D}^{\mathbb{Q}})^{-1}\right]^{-1}, \mathbf{v}_0 + \mathbf{v}_1\right)$,
 - (j) $\boldsymbol{\Sigma}^{-(k)} | \mathbf{U}, \mathbf{b}; \boldsymbol{\mu}^{(k)}, \mathbb{Q}^{-(k)}; \mathbf{v}_0 \stackrel{(52)}{\sim} W_{d_R}\left(\tilde{\mathbb{Q}}^{(k)}, n^k(\mathbf{U}) + \mathbf{v}_0\right)$,
 - (k) $\mathbf{b}_i | \mathbb{Y}_i^N, \mathcal{C}_i; \mathbb{Y}_i^{*,\text{OB}}, U_i; \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1} \stackrel{(54)}{\sim} N_{d_R}\left(\tilde{\mathbf{b}}_i, \left[\tilde{\mathbb{Z}}_i^\top \tilde{\mathbb{Z}}_i + \boldsymbol{\Sigma}^{-(k)}\right]^{-1}\right)$,

where in each step we use the latest generated values of parameters, i.e. either elements of previously generated $\boldsymbol{\Psi}_{[m-1]}$ or currently generated $\boldsymbol{\Psi}_{[m]}$.

3. If $m < M$ then increase $m := m + 1$ and repeat step 2., otherwise end the algorithm.
4. Output: M generated states $\boldsymbol{\Psi}_{[1]}, \dots, \boldsymbol{\Psi}_{[M]}$.

Generated values $\boldsymbol{\Psi}_{[1]}, \dots, \boldsymbol{\Psi}_{[M]}$ are then used for estimation of marginal distributions of individual parameters. Usual methods as cutting first B values (burn-in period) that serve to get to posterior distribution and thinning the rest of the chain to lower the autocorrelation were applied.

For example, the final classification of subject i into one of K latent groups can be based on generated values $U_{i,[1]}, \dots, U_{i,[M]}$. When cutting the burn-in period (of length $B < M$) and thinning (of order t) are applied then subject i can be classified into the most common value among the rest of generated U_i :

$$\hat{U}_i = \text{mod} \left\{ U_{i,[B+m*t+1]}, m \in \left\{ 0, 1, \dots, \left\lfloor \frac{M-B-1}{t} \right\rfloor \right\} \right\}. \quad (27)$$

For practical reasons it also makes sense to actually classify subject into this group only if the corresponding proportion is higher than given threshold, if not, the subject remains unclassified to avoid misclassification.

Another way of classification is presented in the next section, which also covers the situation of newly observed subject.

4.6.3 Label switching

One of the possible disadvantages of MCMC approach is the fact that the latent clusters $1, \dots, K$ do not have fixed meaning. The resulting estimates of class-specific parameters could be arbitrarily permuted while preserving the final model. There is a danger of possible change of the cluster interpretation during a single chain called *label switching*. This problem could be dealt within post-processing phase by following recommendations of Stephens (2000). As we consider low values of K we adopt an algorithm that examines all $K!$ possibilities for each iteration and chooses the one minimizing a given criterion. However, during our experiments no need of permutation has occurred. This may be given by the complexity of our model that prevents any switching of the interpretation during our Gibbs sampling procedure.

4.7 Classification probabilities calculation

Let us consider subject with outcomes \mathbb{Y}_{new} and covariate values \mathcal{C}_{new} , our task is to classify this subject into one of K latent groups. In Section ?? we described mixture probability density function (8) and (??) of all outcomes given covariates for one subject. Derived classification probability $u_{i,k}$ is a function that depends on parameters of interest θ but is free of latent parameters such as $\mathbb{Y}_{\text{new}}^{*,\text{OB}}$, \mathbf{b}_{new} or $\tau_{\text{R}}, \mathbb{Q}^{-1}$. Such a parametric function can be evaluated for each generated state $\theta_{[m]}$ to obtain a chain of posterior classification probabilities. Unfortunately, each evaluation of such function requires integration with respect to random effects \mathbf{b}_{new} and latent numeric outcomes $\mathbb{Y}_{\text{new}}^{*,\text{OB}}$ described in (8):

$$\begin{aligned} & \int \int p\left(\mathbb{Y}_{\text{new}}^{\text{OB}} \mid \mathbb{Y}_{\text{new}}^{*,\text{OB}}; \boldsymbol{\gamma}\right) \cdot p\left(\mathbb{Y}_{\text{new}}^{\text{N}}, \mathbb{Y}_{\text{new}}^{*,\text{OB}} \mid \mathcal{C}_{\text{new}}, \mathbf{b}_{\text{new}}; \boldsymbol{\beta}, \boldsymbol{\tau}\right) \cdot p\left(\mathbf{b}_{\text{new}} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}\right) d\mathbf{b}_{\text{new}} d\mathbb{Y}_{\text{new}}^{*,\text{OB}} = \\ & = \int p\left(\mathbb{Y}_{\text{new}}^{\text{OB}} \mid \mathbb{Y}_{\text{new}}^{*,\text{OB}}; \boldsymbol{\gamma}\right) \cdot \underbrace{\left[\int p\left(\mathbb{Y}_{\text{new}}^{\text{N}}, \mathbb{Y}_{\text{new}}^{*,\text{OB}} \mid \mathcal{C}_{\text{new}}, \mathbf{b}_{\text{new}}; \boldsymbol{\beta}, \boldsymbol{\tau}\right) \cdot p\left(\mathbf{b}_{\text{new}} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}\right) d\mathbf{b}_{\text{new}} \right]}_{p\left(\mathbb{Y}_{\text{new}}^{\text{N}}, \mathbb{Y}_{\text{new}}^{*,\text{OB}} \mid \mathcal{C}_{\text{new}}; \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}\right)} d\mathbb{Y}_{\text{new}}^{*,\text{OB}} \quad (28) \end{aligned}$$

4.7.1 Integration with respect to random effects \mathbf{b}_{new}

Let us first perform the integration of random effects \mathbf{b}_{new} from (11) to obtain distribution of numeric variables unconditioned by random effects. This issue is commonly faced and solved in LME models by realizing that under normality assumption of both numeric outcomes and random effects the unconditioned distribution of outcomes is also normal. Let us stack all outcomes into one long vector:

$$\mathbf{Y}_{\text{new}} = \begin{pmatrix} \vdots \\ \mathbf{Y}_{\text{new}}^{\text{rN}} \\ \vdots \\ \mathbf{Y}_{\text{new}}^{*,\text{rOB}} \\ \vdots \end{pmatrix}, \quad \begin{array}{l} r^{\text{N}} \in \mathcal{R}^{\text{Num}}, \\ r^{\text{OB}} \in \mathcal{R}^{\text{OB}} \end{array} \quad \text{and stack regressors into diagonal matrices}$$

$$\mathbf{X}_{\text{new}} = \begin{pmatrix} \ddots & & & & \\ & \mathbf{X}_{\text{new}}^{\text{rN}} & & & \\ & & \ddots & & \\ & & & \mathbf{X}_{\text{new}}^{\text{rOB}} & \\ & & & & \ddots \end{pmatrix}, \quad \mathbf{Z}_{\text{new}} = \begin{pmatrix} \ddots & & & & \\ & \mathbf{Z}_{\text{new}}^{\text{rN}} & & & \\ & & \ddots & & \\ & & & \mathbf{Z}_{\text{new}}^{\text{rOB}} & \\ & & & & \ddots \end{pmatrix}, \quad \begin{array}{l} r^{\text{N}} \in \mathcal{R}^{\text{Num}}, \\ r^{\text{OB}} \in \mathcal{R}^{\text{OB}}. \end{array}$$

In this section let $\boldsymbol{\beta}, \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}$ and $\boldsymbol{\tau}$ stand for only one representative of the parameter among K groups (similarly as in opening Section ??). Moreover, let $\boldsymbol{\beta}$ and \mathbf{b}_{new} be stacked vectors of fixed and random effects of all outcomes,

i. e.

$$\boldsymbol{\beta} = \begin{pmatrix} \vdots \\ \boldsymbol{\beta}^{r^N} \\ \vdots \\ \boldsymbol{\beta}^{r^{OB}} \\ \vdots \end{pmatrix}, \quad \mathbf{b}_{\text{new}} = \begin{pmatrix} \vdots \\ \mathbf{b}_{\text{new}}^{r^N} \\ \vdots \\ \mathbf{b}_{\text{new}}^{r^{OB}} \\ \vdots \end{pmatrix}, \quad \begin{array}{l} r^N \in \mathcal{R}^{\text{Num}}, \\ r^{OB} \in \mathcal{R}^{OB}. \end{array}$$

Adding dimensions $d = n_{\text{new}} |\mathcal{R}|$, $d^N = n_{\text{new}} |\mathcal{R}^{\text{Num}}|$, $d^{OB} = n_{\text{new}} |\mathcal{R}^{OB}|$

$$\mathbf{Y}_{\text{new}} | \mathbf{b}_{\text{new}} \sim N_d \left(\mathbb{X}_{\text{new}} \boldsymbol{\beta} + \mathbb{Z}_{\text{new}} \mathbf{b}_{\text{new}}, \mathbb{T} = \text{diag} \left\{ \tau_{r^N}^{-1} \mathbb{I}_{n_{\text{new}}}, r^N \in \mathcal{R}^{\text{Num}}, \mathbb{I}_{d^{OB}} \right\} \right)$$

Now we can use notoriously known formulas to obtain unconditioned mean and variance matrix:

$$\begin{aligned} \mathbb{E} \mathbf{Y}_{\text{new}} &= \mathbb{E} (\mathbb{E} [\mathbf{Y}_{\text{new}} | \mathbf{b}_{\text{new}}]) = \mathbb{X}_{\text{new}} \boldsymbol{\beta} + \mathbb{Z}_{\text{new}} \boldsymbol{\mu}, \\ \text{var} \mathbf{Y}_{\text{new}} &= \mathbb{E} (\text{var} [\mathbf{Y}_{\text{new}} | \mathbf{b}_{\text{new}}]) + \text{var} (\mathbb{E} [\mathbf{Y}_{\text{new}} | \mathbf{b}_{\text{new}}]) = \mathbb{T} + \mathbb{Z}^\top \boldsymbol{\Sigma} \mathbb{Z} =: \mathbb{V}, \end{aligned}$$

which results in

$$\mathbf{Y}_{\text{new}} \sim N_d \left(\mathbb{X}_{\text{new}} \boldsymbol{\beta} + \mathbb{Z}_{\text{new}} \boldsymbol{\mu}, \mathbb{T} + \mathbb{Z}^\top \boldsymbol{\Sigma} \mathbb{Z} \right), \quad (29)$$

which can be distribution with completely general covariance structure due to our generality of $\boldsymbol{\Sigma}$. Thus, dependencies among outcomes are captured by this model.

4.7.2 Integration with respect to latent numeric outcomes $\mathbb{Y}_{\text{new}}^{*,OB}$

It remains to perform the following integration:

$$\int p \left(\mathbb{Y}_{\text{new}}^{OB} | \mathbb{Y}_{\text{new}}^{*,OB}; \boldsymbol{\gamma} \right) \cdot p \left(\mathbb{Y}_{\text{new}}^N, \mathbb{Y}_{\text{new}}^{*,OB} | \mathcal{C}_{\text{new}}; \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1} \right) d\mathbb{Y}_{\text{new}}^{*,OB},$$

which is in fact integration of multivariate normal density within the bounds given by thresholds $\boldsymbol{\gamma}$ and observed ordinal and binary outcomes. Only the truly numerical part survives this integration intact, therefore, we aim to separate it from the rest.

Let us divide vector \mathbf{Y}_{new} on the truly numerical part and the vector of latent numeric variables for ordinal and binary outcomes:

$$\begin{pmatrix} \mathbf{Y}_{\text{new}}^N \\ \mathbf{Y}_{\text{new}}^{*,OB} \end{pmatrix} \sim N_{d^N + d^{OB}} \left(\begin{pmatrix} \boldsymbol{\eta}^N \\ \boldsymbol{\eta}^{OB} \end{pmatrix} = \begin{pmatrix} \mathbb{X}_{\text{new}}^N \boldsymbol{\beta}^N + \mathbb{Z}_{\text{new}}^N \boldsymbol{\mu}^N \\ \mathbb{X}_{\text{new}}^{OB} \boldsymbol{\beta}^{OB} + \mathbb{Z}_{\text{new}}^{OB} \boldsymbol{\mu}^{OB} \end{pmatrix}, \begin{pmatrix} \mathbb{V}_N & \mathbb{V}_{NOB} \\ \mathbb{V}_{OBN} & \mathbb{V}_{OB} \end{pmatrix} \right). \quad (30)$$

Since the latent part is about to be integrated out we exploit properties of multivariate normal distribution to obtain

$$\begin{aligned} \mathbf{Y}_{\text{new}}^N &\sim N_{d^N} \left(\boldsymbol{\eta}^N, \mathbb{V}_N \right), \\ \mathbf{Y}_{\text{new}}^{*,OB} | \mathbf{Y}_{\text{new}}^N = \mathbf{y}_{\text{new}}^N &\sim N_{d^{OB}} \left(\boldsymbol{\eta}'_{OB}, \mathbb{V}'_{OB} \right), \end{aligned}$$

where $\boldsymbol{\eta}'_{OB} = \boldsymbol{\eta}^{OB} + \mathbb{V}_{OBN} \mathbb{V}_N^{-1} (\mathbf{y}_{\text{new}}^N - \boldsymbol{\eta}^N)$ and $\mathbb{V}'_{OB} = \mathbb{V}_{OB} - \mathbb{V}_{OBN} \mathbb{V}_N^{-1} \mathbb{V}_{NOB}$ and corresponding probability distribution function decomposes into the product

$$\begin{aligned} p \left(\mathbf{Y}_{\text{new}}^N, \mathbf{Y}_{\text{new}}^{*,OB} | \mathcal{C}_{\text{new}}; \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1} \right) &= \\ &= p \left(\mathbf{Y}_{\text{new}}^{*,OB} | \mathbf{Y}_{\text{new}}^N, \mathcal{C}_{\text{new}}; \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1} \right) \cdot p \left(\mathbf{Y}_{\text{new}}^N | \mathcal{C}_{\text{new}}; \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1} \right), \end{aligned}$$

second part of which is invariant of $\mathbf{Y}_{\text{new}}^{*,OB}$ and, therefore, avoids integration:

$$p \left(\mathbf{Y}_{\text{new}}^N | \mathcal{C}_{\text{new}}; \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1} \right) \cdot \underbrace{\int p \left(\mathbf{Y}_{\text{new}}^{OB} | \mathbf{Y}_{\text{new}}^{*,OB}; \boldsymbol{\gamma} \right)}_{\text{thresholding (3)}} \cdot \underbrace{p \left(\mathbf{Y}_{\text{new}}^{*,OB} | \mathbf{Y}_{\text{new}}^N, \mathcal{C}_{\text{new}}; \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1} \right)}_{\text{pdf of } N_{d^{OB}}(\boldsymbol{\eta}'_{OB}, \mathbb{V}'_{OB})} d\mathbf{Y}_{\text{new}}^{*,OB}.$$

I zde je třeba se rozhodnout, jak s tím tlustým \mathbf{Y} či \mathbb{Y} .

It remains to integrate the product of two functions, first one of which just declares lower and upper integration bounds and the second one is the probability density function of multivariate normal distribution $N_{d^{\text{OB}}}(\boldsymbol{\eta}'_{\text{OB}}, \mathbb{V}'_{\text{OB}})$. For each individual categorical outcome $r^{\text{OB}} \in \mathcal{R}^{\text{OB}}$ and observation $j \in \{1, \dots, n_{\text{new}}\}$ the value $y_{\text{new},j}^{\text{OB}} = l$ determines an interval given by the corresponding pair of γ parameters, see (3):

$$Y_{\text{new},j}^{\text{OB}} = l \implies Y_{\text{new},j}^{*,r^{\text{OB}}} \in \left(\gamma_l^{\text{OB}}, \gamma_{l+1}^{\text{OB}} \right] =: \left(e_j^{\text{OB}}, f_j^{\text{OB}} \right].$$

If we denote the resulting Cartesian product of these intervals as $\square(\mathbf{Y}_{\text{new}}^{\text{OB}}) \subset \mathbb{R}^{d^{\text{OB}}}$, i.e.

$$\square(\mathbf{Y}_{\text{new}}^{\text{OB}}) = (\mathbf{e}, \mathbf{f}] = \bigotimes_{r^{\text{OB}} \in \mathcal{R}^{\text{OB}}} \bigotimes_{j=1}^{n_{\text{new}}} \left(e_j^{\text{OB}}, f_j^{\text{OB}} \right] = \bigotimes_{r^{\text{OB}} \in \mathcal{R}^{\text{OB}}} \bigotimes_{j=1}^{n_{\text{new}}} \left(\gamma_{Y_{\text{new},j}^{*,r^{\text{OB}}}}^{\text{OB}}, \gamma_{Y_{\text{new},j}^{*,r^{\text{OB}}}+1}^{\text{OB}} \right],$$

then the remaining integral can be written in the form

$$I(\square) = \int_{\square(\mathbf{Y}_{\text{new}}^{\text{OB}})} p\left(\mathbf{Y}_{\text{new}}^{*,\text{OB}} \mid \mathbf{Y}_{\text{new}}^{\text{N}}, \mathcal{C}_{\text{new}}; \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}\right) d\mathbf{Y}_{\text{new}}^{*,\text{OB}} = \int_{\mathbf{e}}^{\mathbf{f}} p(\mathbf{y} \mid \boldsymbol{\eta}'_{\text{OB}}, \mathbb{V}'_{\text{OB}}) d\mathbf{y}, \quad (31)$$

where $p(\mathbf{y} \mid \boldsymbol{\eta}'_{\text{OB}}, \mathbb{V}'_{\text{OB}})$ denotes the probability density function of multivariate normal distribution $N_{d^{\text{OB}}}(\boldsymbol{\eta}'_{\text{OB}}, \mathbb{V}'_{\text{OB}})$.

In order to compute such integral we adopted an effective algorithm presented by Genz (1992) which is also based on MCMC sampling. Since the approximation of such an integral is needed K -times for each generated state of our Gibbs sampling, the procedure is considerably time-consuming. We exploit the implemented function `pvmnorm` in \mathbb{R} package `mvtnorm` (Genz et al, 2019).

Altogether, the probability of being classified into group k is proportional to

$$\begin{aligned} u_{\text{new},k}(\boldsymbol{\theta}) &\propto w_k \cdot p\left(\mathbf{Y}_{\text{new}}^{\text{N}} \mid \mathcal{C}_{\text{new}}; \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{-(k)}\right) \cdot I_k(\square), \\ &= \frac{w_k \cdot p\left(\mathbf{Y}_{\text{new}}^{\text{N}} \mid \mathcal{C}_{\text{new}}; \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{-(k)}\right) \cdot I_k(\square)}{\sum_{k'=1}^K w_{k'} \cdot p\left(\mathbf{Y}_{\text{new}}^{\text{N}} \mid \mathcal{C}_{\text{new}}; \boldsymbol{\beta}^{(k')}, \boldsymbol{\tau}^{(k')}, \boldsymbol{\mu}^{(k')}, \boldsymbol{\Sigma}^{-(k')}\right) \cdot I_{k'}(\square)} \end{aligned} \quad (32)$$

where I_k is the integral (13) computed with set of parameters $\boldsymbol{\theta}^{(k)}$ belonging to group k . These classification probabilities are then computed for all (or just a subset due to burn-in period and thinning) generated states $\boldsymbol{\theta}_{[1]}, \dots, \boldsymbol{\theta}_{[M]}$ and can be used to describe the posterior distribution. Subjects then can be classified to the class with convincingly high value of posterior mean estimator $\hat{u}_{\text{new},k}$. The subject might remain unclassified, if neither of the classes seems to dominate others.

4.7.3 Deviance

Throughout the paper we treated the number of latent classes K fixed as to be selected by statistician in advance. In most circumstances, however, there is no prior knowledge of the suitable value of K to be used. Usual practise in this situation is to try several values of K and choose the one that optimizes one of the known criteria that measure the goodness of fit to the data and penalizes the number of estimated unknown parameters $\boldsymbol{\theta}$. We present here a deviance criterion that is solely based on the goodness of fit measured by the log-likelihood function.

Deviance is defined as

$$D(\boldsymbol{\theta}; \mathbb{Y}, \mathcal{C}) = -2 \log p(\mathbb{Y} \mid \mathcal{C}; \boldsymbol{\theta}) = -2 \sum_{i=1}^n \log p(\mathbb{Y}_i \mid \mathcal{C}_i; \boldsymbol{\theta}). \quad (33)$$

Contribution of individual i to the deviance can be expressed as

$$-2 \log p(\mathbb{Y}_i \mid \mathcal{C}_i; \boldsymbol{\theta}) = -2 \log \left[\sum_{k=1}^K \int \int p\left(\mathbb{Y}_i, \mathbb{Y}_i^{*,\text{OB}}, \mathbf{b}_i, U_i = k \mid \mathcal{C}_i; \boldsymbol{\theta}\right) d\mathbf{b}_i d\mathbb{Y}_i^{*,\text{OB}} \right] \quad (34)$$

which includes the integration (11) which has been described above. Once we use analogous notation we can simply write

$$D(\boldsymbol{\theta}; \mathbb{Y}, \mathcal{C}) = -2 \sum_{i=1}^n \log \left[\sum_{k=1}^K w_k \cdot p \left(\mathbb{Y}_i^{\mathbb{N}} \mid \mathcal{C}_i; \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{-(k)} \right) \cdot I_k \left(\square \left(\mathbb{Y}_i^{\text{OB}} \right) \right) \right] \quad (35)$$

where the denominator of (14) was inserted into the logarithm. In order to compute the deviance (and all related criteria) we need to basically compute classification probabilities for every individual. Deviance is, as a matter of fact, a parametric function of $\boldsymbol{\theta}$ and has, therefore, its posterior distribution, which can be estimated by calculation of deviance parametric function on a sample $\boldsymbol{\theta}_{[1]}, \dots, \boldsymbol{\theta}_{[M]}$. As the computation of classification probabilities is time-consuming we resort to thinning.

5 Simulation

In this section we design and perform a simulation study in order to demonstrate the functionality of our proposed model. Data consisting of a numeric, a binary and an ordinal variable were generated while assuming different types of random effects structure. The only parameter distinguishing the latent groups ($K = 2$ or $K = 3$) was the parameter connected to the parametrization of time, i.e. intercept or slope. Parameters describing the covariance structure ($\boldsymbol{\tau}$ and $\boldsymbol{\Sigma}^{-1}$) were held equal for all latent groups.

5.1 Simulation design

Each type (numeric, binary and ordinal) is represented by just one longitudinally measured variable ($Y_{i,j}^{\mathbb{N}}, Y_{i,j}^{\mathbb{B}}, Y_{i,j}^{\mathbb{O}}$, $i = 1, \dots, n$ and $j = 1, \dots, n_i$). We set the number of observations per one subject n_i to be fixed at $n_i = 4$ for each of the n subjects, $n \in \{100, 500, 1000\}$, which corresponds to the same amount of observations per household available in the EU-SILC data.

The part of the predictor which is common to all types of variables is of the form

$$1 \cdot X_{i,j}^1 - 2 \cdot X_{i,j}^2, \quad \text{where } X_{i,1}^1 = \dots = X_{i,4}^1 \stackrel{\text{iid}}{\sim} \text{Alt}(0.5) \quad \text{and} \quad X_{i,j}^2 \stackrel{\text{iid}}{\sim} \text{Unif}(0, 1).$$

Then we suppose that each subject has his set of times $0 < t_{i,1} < t_{i,2} < t_{i,3} < t_{i,4} < 1$ which were generated as an ordered sample from uniform distribution over interval $(0, 1)$. We suppose the linear parametrization of time, which is, however, different subject to the structure of random effects. We consider three scenarios

1. (r=intercept) - $b_{0,i} + \beta_1 t_{i,j}$, random intercept term and fixed slope,
2. (r=slope) - $\beta_0 + b_{1,i} t_{i,j}$, fixed intercept term and random slope,
3. (r=both) - $b_{0,i} + b_{1,i} t_{i,j}$, both intercept and slope are random effects.

We keep the same random effects structure for all outcome types. Therefore, the random effects of i -th subject are multivariate normal of dimension 3 (cases 1. and 2.) or 6 (case 3.). Its variance matrix $\boldsymbol{\Sigma}$ was adequately chosen matrix of non-diagonal form, more details can be found in supplement materials.

Another 3 types of scenarios arise from the types of differences among $K = 2$ or $K = 3$ latent groups:

- a) (d=intercept) - only in the intercept term - $\beta_0^{(k)}$ (case 2.) and $\boldsymbol{\mu}_0^{(k)}$ (case 1. and 3.) are class-specific, but slope parameters β_1 and $\boldsymbol{\mu}_1$ are not,
- b) (d=slope) - only in the slope - $\beta_1^{(k)}$ (case 1.) and $\boldsymbol{\mu}_1^{(k)}$ (case 2. and 3.) are class-specific, but intercept terms β_0 and $\boldsymbol{\mu}_0$ are not,
- c) (d=both) - both in the intercept term and the slope - $\beta_0^{(k)}, \beta_1^{(k)}, \boldsymbol{\mu}_0^{(k)}$ and $\boldsymbol{\mu}_1^{(k)}$ are class-specific.

These three types of differences are combined with the three types of random effects structure creating 9 different scenarios which are examined for $K = 2, 3$ and different sample sizes n . The values of intercept and slope for each of the 9 scenarios were chosen in different ways to obtain clusters distinguishable by eye (see Figure 2 for the case $K = 3$), the true values of intercept and slope parameters can be found in Tables S?? and S?? in supplement materials.

The group allocation indicator U_i was always generated from uniform distribution which results in clusters of comparable sizes. All (latent) numeric outcomes were sampled with unit variance $\tau = 1$. Binary variable was obtained by threshold $\gamma_1^{\mathbb{B}} = 1$ and the ordinal one by thresholds $\gamma_1^{\mathbb{O}} = -1$ and $\gamma_2^{\mathbb{O}} = 2$.

Each scenario under given K and n was replicated 200-times to explore the properties of resulting estimators. Estimation is based on $M = 10000$ sampled states from posterior distribution. To spare computation time, thinning of order 10 was applied for classification probabilities leading to 1000 calculated probabilities per subject. We are also interested in overall probability that subject from k -th cluster is correctly classified to cluster k . For this reason we compute for each k arithmetic mean \bar{p}_k of our estimates $\widehat{U}_{i,k}$ of posterior probabilities of belonging to cluster k across all cluster members:

$$\bar{p}_k = \frac{1}{|i : U_i = k|} \sum_{i:U_i=k} \widehat{U}_{i,k}.$$

These probabilities are also estimated dynamically, i.e. using only limited amount of information, see 4.4.

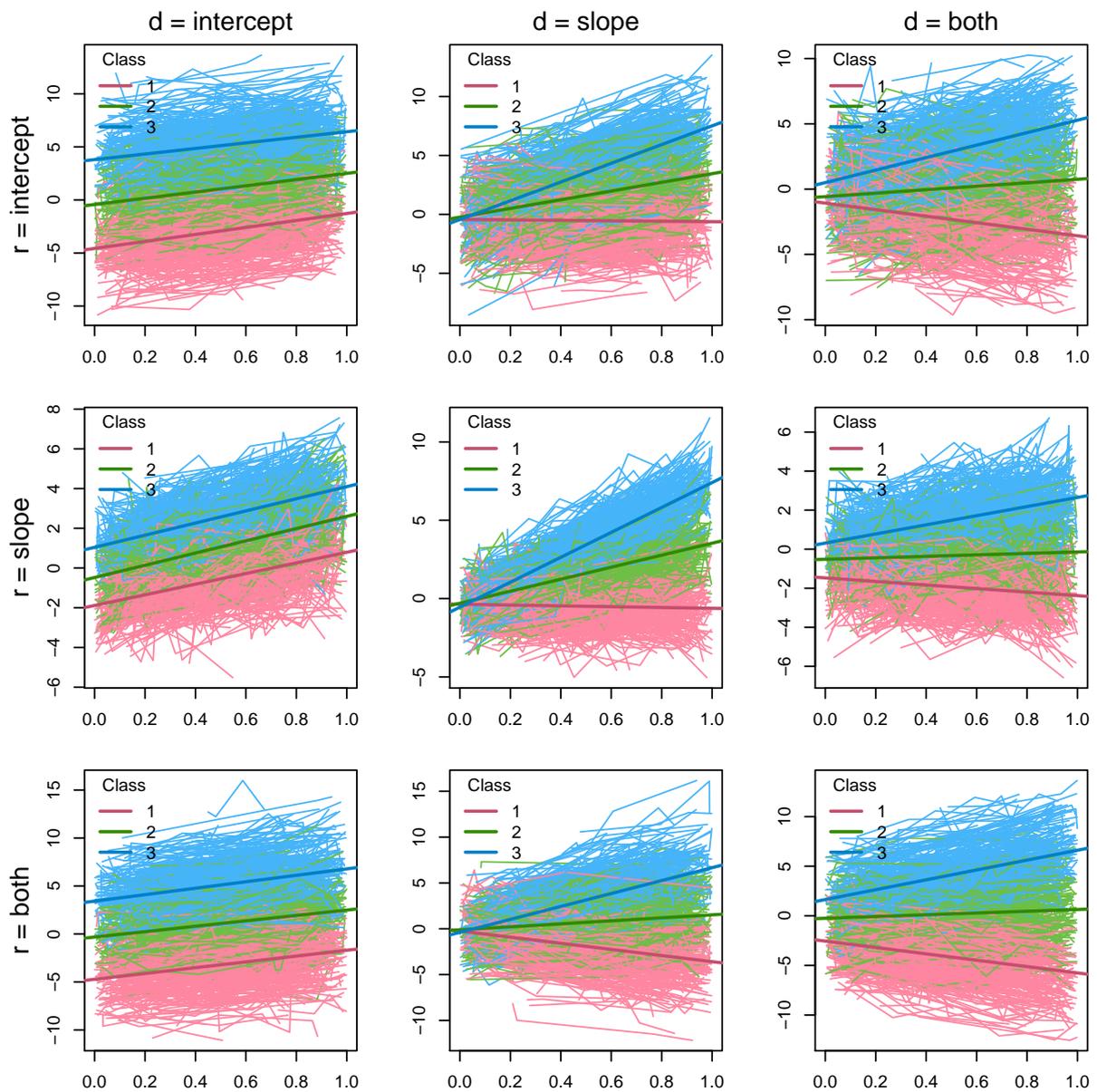


Fig. 2: The samples of a numeric outcome distinguishing different scenario types (row difference, column structure of random effects) when $K = 3$ latent groups are supposed.

5.2 Results

Let us first focus on the estimation properties of the proposed modelling procedure, which is shown by Figure 3. Colours distinguish the estimates in different classes (clusters?) and the corresponding true values of intercept and slope parameters are captured by dashed lines. Grey colour depicts the true value shared by all classes. Each segment represents 2.5% and 97.5% quantiles of 200 times replicated estimator and full circle represents corresponding mean. Figure 3 provides estimates of parameters belonging to ordinal outcome only, plots for numeric and binary are postponed to supplement materials.

Figure 3 demonstrates that the proposed procedure is capable of quite precise estimation of parameters despite the latent modelling and thresholding concept. In most cases it successfully discovers the difference among classes as intervals of different colours tend not to overlap each other. There is apparent a decreasing trend in standard deviation as n increases which is disrupted only when corresponding estimate does not reach the true value. This phenomenon occurs mostly in the estimation of intercept term when it is considered to be random and different among clusters at the same time. Such a behaviour can be spotted also for the class-specific slope term when both intercept and slope term are random effects. In these situations, the estimates are shrunk towards a mean of the true values. This might be a result of combination of incapability of discrimination between classes for low value of n and the fact that LME usually tend to shrink random effects to zero. In the case of $K = 3$, this effect does not fully vanish even for $n = 1000$, see the row *both* and column *intercept*. However, it seems that the large number of subjects n can overcome this issue, which we rely on in the further real data analysis.

As the ability to discover differences among classes has been verified we now proceed to examine the classification abilities. Table 1 contains the percentages of correctly classified subjects (using the HPD interval rule) among all n of them across 200 replications. This percentage differs scenario by scenario as the random structures and differences among classes interact in different ways leading to diverse success rates. For example, the case with class-specific random slope successfully classifies the vast majority of subjects for both $K = 2$ and $K = 3$, which is in agreement with the strict separation in the corresponding plot of Figure 2. Classification does not work satisfactorily in the problematic cases discussed above. Since for the low values of n the difference between classes is not estimated to be as strict as it should be, much larger percentage of subjects is kept unclassified in such cases. With increasing n the percentage of unclassified subjects rapidly decreases and converts mainly into correctly classified category. Under all scenarios we managed to keep the misclassification rate very low, always under 10%. The unclassified proportion is also much higher for $K = 3$ as one of the classes (green) is surrounded from both sides which significantly declines the ability to distinguish among classes, revisit Figure 2 for illustration.

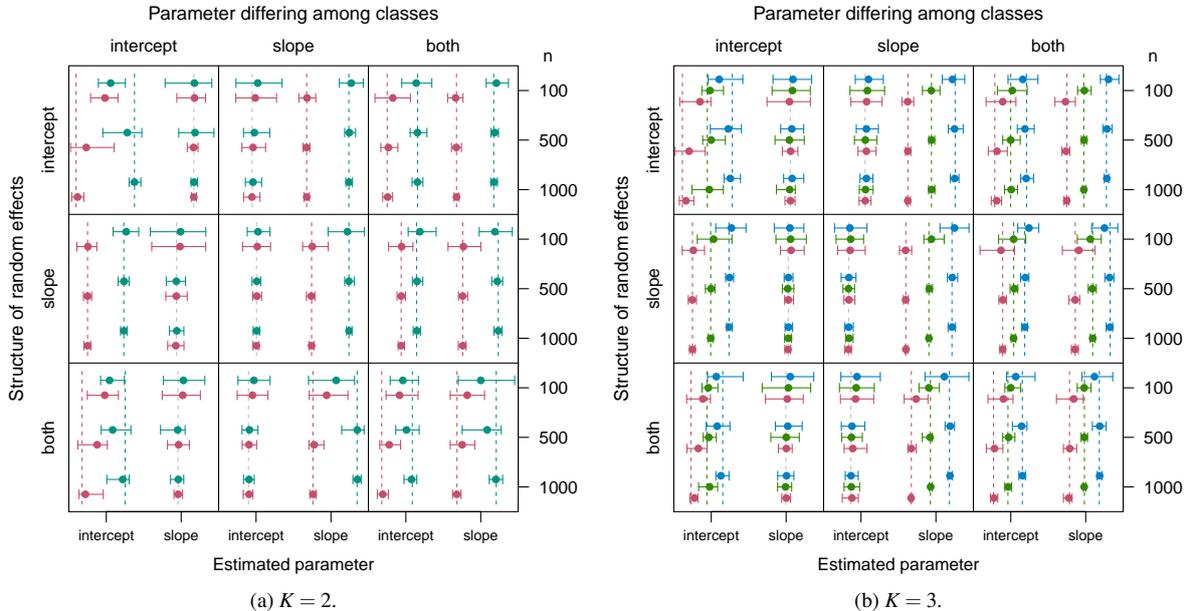


Fig. 3: 95% quantile bounds and means for the intercept and slope parameters for ordinal outcome under different random effects structures and differences between classes. True values of the parameters are depicted by dashed lines (grey if common to all classes).

Table 1: Percentages (standard deviation) of correctly classified, unclassified and misclassified subjects (using the HPD interval rule) for several choices of n , K , structure of random effects and class differences in 200 replications.

| ¹ r | ² d | n | $K = 2$ | | | $K = 3$ | | |
|-------------------|-------------------|------|-------------|-------------|------------|-------------|-------------|------------|
| | | | Correct [%] | Uncl. [%] | Miscl. [%] | Correct [%] | Uncl. [%] | Miscl. [%] |
| intercept | intercept | 100 | 27.0 (17.2) | 63.2 (25.4) | 9.8 (13.7) | 23.0 (17.5) | 70.2 (21.0) | 6.8 (9.4) |
| | | 500 | 62.5 (27.2) | 33.0 (27.3) | 4.4 (3.8) | 44.3 (20.6) | 50.8 (22.1) | 4.9 (4.4) |
| | | 1000 | 85.1 (6.7) | 10.1 (7.1) | 4.8 (0.9) | 58.6 (16.9) | 35.5 (17.2) | 6.0 (3.1) |
| intercept | slope | 100 | 76.8 (5.4) | 20.3 (5.5) | 2.9 (1.9) | 56.0 (8.5) | 40.4 (8.9) | 3.6 (2.4) |
| | | 500 | 86.1 (1.8) | 8.9 (1.8) | 5.0 (1.0) | 74.6 (2.0) | 19.0 (2.1) | 6.4 (1.2) |
| | | 1000 | 87.5 (1.1) | 6.7 (0.9) | 5.9 (0.7) | 78.2 (1.5) | 13.8 (1.5) | 8.0 (0.8) |
| intercept | both | 100 | 86.5 (4.4) | 12.0 (4.4) | 1.5 (1.1) | 58.0 (9.4) | 38.5 (10.2) | 3.4 (2.2) |
| | | 500 | 92.9 (1.4) | 4.5 (1.1) | 2.6 (0.7) | 76.9 (2.5) | 16.5 (2.5) | 6.7 (1.1) |
| | | 1000 | 93.8 (0.8) | 3.3 (0.6) | 2.9 (0.5) | 79.4 (1.6) | 12.8 (1.6) | 7.8 (0.8) |
| slope | intercept | 100 | 96.2 (2.6) | 3.4 (2.5) | 0.4 (0.6) | 61.2 (15.5) | 36.4 (15.7) | 2.3 (1.8) |
| | | 500 | 97.9 (0.5) | 1.5 (0.5) | 0.6 (0.4) | 87.6 (2.2) | 9.2 (2.2) | 3.2 (0.7) |
| | | 1000 | 98.3 (0.4) | 0.9 (0.3) | 0.8 (0.3) | 90.2 (1.2) | 6.2 (1.1) | 3.6 (0.5) |
| slope | slope | 100 | 80.1 (20.4) | 16.3 (19.0) | 3.6 (8.7) | 85.7 (13.5) | 13.3 (13.6) | 1.0 (1.2) |
| | | 500 | 92.8 (1.5) | 4.6 (1.4) | 2.6 (0.7) | 94.9 (1.2) | 3.6 (1.0) | 1.5 (0.5) |
| | | 1000 | 93.9 (0.9) | 3.3 (0.7) | 2.8 (0.5) | 95.5 (0.7) | 2.6 (0.5) | 1.9 (0.4) |
| slope | both | 100 | 85.3 (18.0) | 13.8 (18.0) | 0.9 (0.9) | 62.2 (23.5) | 35.8 (23.4) | 2.0 (2.7) |
| | | 500 | 96.2 (1.0) | 2.6 (0.9) | 1.3 (0.6) | 92.4 (1.7) | 5.5 (1.5) | 2.1 (0.8) |
| | | 1000 | 96.7 (0.6) | 1.8 (0.4) | 1.5 (0.4) | 93.3 (0.9) | 4.1 (0.9) | 2.5 (0.5) |
| both | intercept | 100 | 18.8 (13.7) | 76.0 (16.6) | 5.2 (7.2) | 18.7 (15.2) | 78.1 (16.7) | 3.2 (4.1) |
| | | 500 | 35.4 (25.2) | 58.7 (27.2) | 6.0 (8.5) | 30.6 (18.5) | 65.1 (20.5) | 4.3 (3.9) |
| | | 1000 | 70.5 (22.4) | 24.3 (23.4) | 5.2 (1.9) | 46.4 (12.1) | 48.2 (13.9) | 5.4 (2.4) |
| both | slope | 100 | 16.2 (13.2) | 79.2 (16.8) | 4.5 (6.0) | 23.4 (22.3) | 74.9 (23.6) | 1.6 (2.3) |
| | | 500 | 69.7 (18.1) | 24.7 (19.4) | 5.6 (2.0) | 69.8 (13.4) | 25.2 (14.4) | 5.0 (1.4) |
| | | 1000 | 80.5 (3.0) | 12.0 (3.3) | 7.4 (1.2) | 81.1 (2.2) | 11.9 (2.1) | 7.0 (0.8) |
| both | both | 100 | 16.7 (14.5) | 80.3 (17.3) | 3.0 (5.5) | 19.4 (19.8) | 79.7 (20.7) | 0.9 (1.4) |
| | | 500 | 43.6 (30.5) | 53.3 (32.3) | 3.0 (2.8) | 66.3 (19.6) | 29.1 (21.1) | 4.5 (1.9) |
| | | 1000 | 80.3 (10.5) | 13.5 (11.1) | 6.2 (1.2) | 80.9 (3.3) | 12.1 (3.5) | 7.0 (1.0) |

¹ Structure of random effects.

² Difference among classes.

Now we discuss the classification properties when limited amount of information is given. Dynamically calculated probabilities, i.e. using just first $j \in \{1, \dots, n_i\}$ observations (see more in 4.4), are shown in Figure 4. It shows mean and quantile bounds of dynamically calculated mean probabilities \bar{p}_2 based on 200 replications of experiments with $K = 3$ clusters. Class 2 has been chosen for demonstration as it is the middle one that overlaps the other two, which covers the most problematic case (with respect to successful classification). The other choices of k and K (with much higher probabilities) can be found in the supplement. If a difference among classes lies in the random intercept term only, then there seems to be no improvement with any additional observation. However,

in other scenarios the probability improves with any additional observation from later times as they help to fit the corresponding medium slope value better. That results in rejecting the low and the extremely large slope values of other classes, and therefore increasing the probability of classification towards the true middle class. It also improves with increasing number of subjects n since the classes are then better distinguished.

The simulation study was conducted on cluster consisting of CPU units: Intel(R) Xeon(R) CPU E5-2620 v2 @ 2.10GHz, 64GB RAM. The average computation time of generating a chain of $M = 10000$ followed by much more demanding computation of 1000 classification probabilities for all n subjects could not fit under an hour even for the lowest values of $n = 100$ and $K = 2$ (around 80 minutes). The number of calls of `pmvnorm` seem to influence the computational time the most. As expected, in the case of $n = 500$ or $n = 1000$ the procedure requires 5 or 10-times more time than if $n = 100$. Similarly, addition of another cluster increases the computational time as well due to one additional call of `pmvnorm` function per subject and set of θ parameters to calculate corresponding classification probabilities. The most challenging combination of $n = 1000$ and $K = 3$ took around 1200 minutes (20 hours).

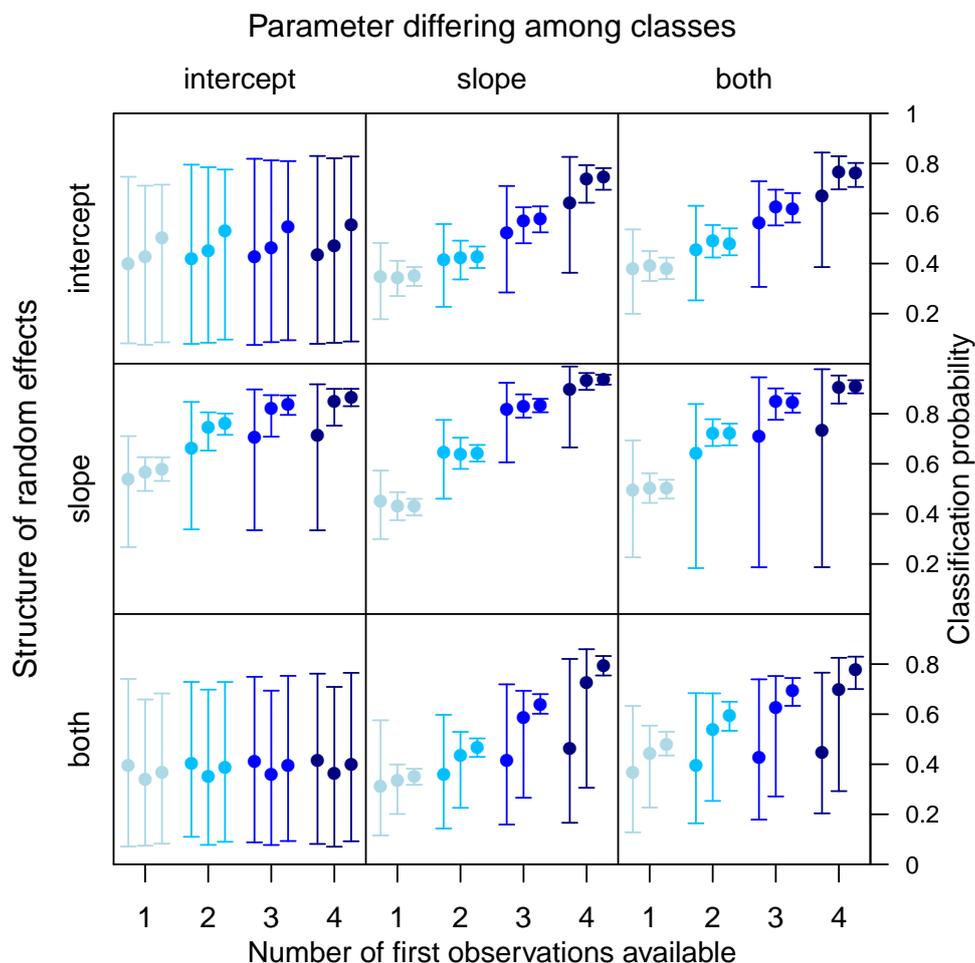


Fig. 4: Subjects of class 2 when $K = 3$. The mean and 2.5% and 97.5% quantile of mean classification probabilities \bar{p}_2 towards the true class calculated dynamically using only first $j \in \{1, 2, 3, 4\}$ observations under several random effects structure and difference among $K = 3$ classes. Three lines of the same colour in one cell correspond to increasing values of $n \in \{100, 500, 1000\}$.

6 Application to EU-SILC data

We will now apply the proposed methodology to find temporal patterns for households in the EU-SILC database in the Czech Republic as outlined in the introduction. The chosen time period 2005 – 2016 covers the economical crisis and we expect it to heavily impact the budget of households leading to different ways of coping with the crisis. From those who were not affected and prosper on, to those who suffer unpleasant consequences. We chose one numeric, binary and ordinal variables that reflect the financial situation of a household the most and we aim to discover several different patterns in their evolution while modelling them jointly.

6.1 Data description

First, we need to delve into the data gathering mechanism which is crucial for appropriate interpretation of results. The EU-SILC longitudinal study follows a rotational design - rotating a part of the sample from one year to the next and retaining the other part unchanged. The study in the Czech Republic was launched in 2005 with more than 7 000 households. Each following year about a quarter of households in the study were dropped and replaced by newly entering ones. Apart from the natural leave from the study, households were followed no longer than 4 years. Since the primary focus is on the evolution part we use for the analysis only the households that were interviewed indeed $n_i = 4$ times. This decision reduces the number of total households used for the analysis to $n = 20323$.

The analysis will be performed on the following outcomes:

- i) Total disposable income (numeric),
- ii) Capacity to afford paying for one week annual holiday away from home (binary - yes/no),
- iii) Financial burden of the total housing cost (ordinal - a heavy burden/a slight burden/no burden at all).

All-year income (in EUR) of the household follows heavily skewed distribution. Therefore, we rather work with its logarithmic transformation which suits our LME assumptions much better. Binary outcome refers to the affordability and to the actual meaning ‘ability to pay’ regardless of whether the household wants it. Ordinal outcome was filled subjectively by the respondent to assess his/her feeling about the extent to which housing costs are a financial burden to the household. For obvious reasons these three cannot be considered as completely independent random variables.

Data contains information about the year and month of the interview (CZE data keep only the quarter - either Q1 or Q2). We construct the time variable as the number of years past the beginning of 2005 which limits the time into the interval $[0, 12)$. For the regression part of the model we will also use the weighted family size which is a sum of weights of each of the household members. Adult person in the role of the head of the family has weight 1 and others have either weight 0.5 or 0.3 depending on whether they are older or younger than 14, respectively.

6.2 Model structure

Since the three outcomes are strongly related we cannot model them separately using independent models. Our proposed joint model capable of capturing the relationships is more than suitable for modelling of these outcomes.

We suppose the LME model of the same structure of both fixed and random component for numeric outcome and latent numeric counterparts of binary and ordinal outcomes. Being aware of possible change in the evolution of these outcomes within the time period 2005 – 2016 we decided to parametrize the effect of time by B-spline of order 3 with knots in the years 2005, 2008, 2010 and 2017, which leads to six β parameters including the intercept. This fixed part of the model is extended by the weighted family size as an additional regressor. The random effects structure, which is also responsible for the covariance structure among outcomes, is simply composed of the zero mean random intercept term which allows households to evolve on different level than others. The model formula for j -th observation of i -th household at time $t_{i,j}$ is then:

$$\underbrace{\beta_0^r + \beta_1^r B_1(t_{i,j}) + \dots + \beta_5^r B_5(t_{i,j}) + \beta_6^r S_{i,j}}_{\text{fixed effects}} + \underbrace{b_{0,i}^r}_{\text{random effects}}, \quad r \in \{N, B, O\},$$

where B_1, \dots, B_5 are B-spline functions corresponding to spline of order 3 with knots at 0, 3, 5 and 12 that does not include the intercept, $S_{i,j}$ is the weighted family size and $\mathbf{b}_{0,i} = (b_{0,i}^N, b_{0,i}^B, b_{0,i}^O)^\top$ is the three-dimensional mean-zero vector of random intercepts.

We suppose that the hidden latent groups of households with different evolution patterns share the same covariance properties, meaning that parameters Σ and τ are supposed to be common for all classes. Hence, the relationships among outcomes are kept the same in all classes. However, we suppose them to differ in the $\beta^{(k)}$ coefficients, which mainly describes the different evolution patterns captured by splines. The revealed clusters will then be characterized according to the shape and level of resulting curve.

6.3 Results

Contrary to the previous simulation study we had to be more careful when sampling from the posterior distribution. Initial values of all the unknown parameters and latent variables were randomly generated to obtain different starting points for sampled chains. The progress in each of the model parameters was visually monitored every 10000 steps in order to suggest a reasonable choice of initial values for the subsequent continuation of sampling. Chains required up to hundred thousands iterations until the visual stationarity in all of the aspects was reached. The slow convergence was mainly caused by the threshold parameters γ due to almost negligible steps. A final chain of length $M = 10000$ used for the analysis and results interpretation was sampled only after such visual stationarity was verified and then checked for label-switching issues. In the calculation of classification probabilities and deviance we again thinned the chain by 10 to spare the computation time.

Following Section 4.7.3, we applied this methodology under several different choices of the number of hidden clusters K and examined the posterior distribution of deviance in hope of selecting the most suitable one. Ideally, we would search for the K such that the decrease of deviance becomes negligible. Although, some improvement in decrement of deviance is visible in Figure 5, we can also notice that the solution for $K = 4$ surprisingly achieved lower deviance than the one for $K = 5$. Higher choice of K led to even lower deviance, however, first signs of overfitting appeared with it. Small groups of households of extraordinary and very specific behaviour emerged.

The solution for $K = 5$ already contains an example of such rare clusters. Figure 6 contains estimated spline curves for logarithm of total disposable income for each of the considered solutions. Case $K = 1$ shows us a general increasing trend flattening after the year 2011 (red curve) that seems to be followed by majority of households even for higher K . With $K = 3$ there appears new violet cluster that follows the same shape of the general curve but on much higher level. Hence, it represents about 5–9 % of households having high income at their disposal. The solution for $K = 2$ actually started with parallel curves of the same shape, however, it slowly transformed one of the clusters into a very rare cluster of U-shaped trend (blue curve). For $K \geq 5$ such cluster appears again accompanied with a golden cluster that behaves reversely. This is why this solution should be rather viewed as an extension of $K = 3$ solution. However, situation $K = 4$ avoids these clusters of extreme behaviour and additionally covers a turquoise cluster representing more than 10 % of households of very low disposable income. This cluster seem to be the reason why this solution defeats $K = 5$ in terms of deviance. Unfortunately, it brings more confusion

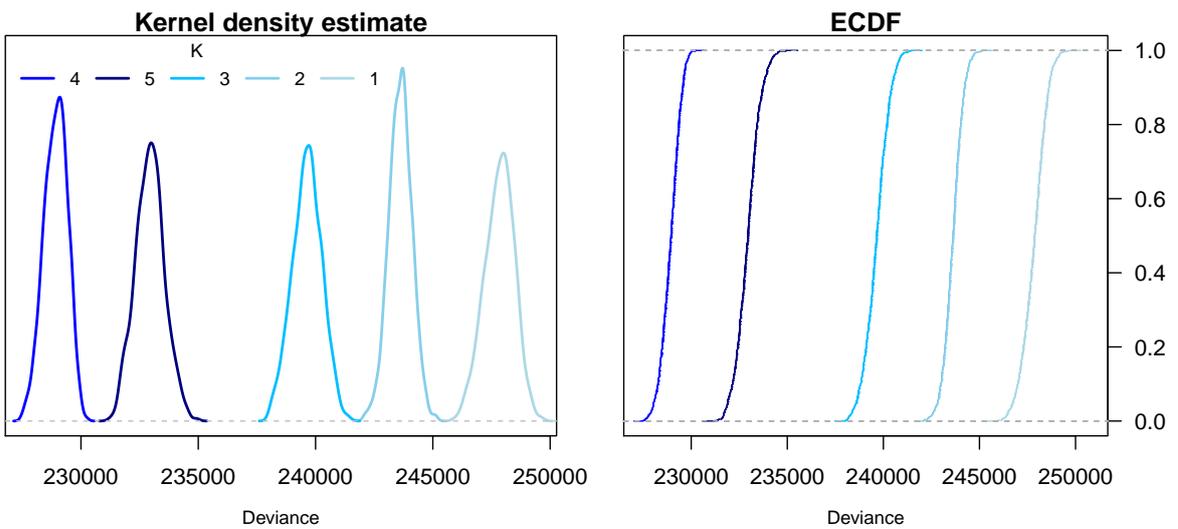


Fig. 5: Comparison of posterior distribution of deviances based on model with number of clusters $K = 1, 2, 3, 4, 5$.

in terms of classification (25 % of households remained unclassified) as the red, green and turquoise clusters do not substantially differ. To interpret such clusters more precisely we should not forget the other used outcomes. The resulting spline shapes for the affordability of week holiday and the financial burden of total housing cost can be found in the supplement, (Figures S12–S13).

Now we discuss the case $K = 5$ in more detail. As our goal is to find different patterns in evolution we may actually be even interested in blue and golden cluster of extreme antagonistic behaviour. These clusters cover households undergoing some substantial transformation which may indeed be what we aim to identify. However, we must not forget the fact that households were followed only for 4 consecutive years. Therefore, the blue cluster should be interpreted rather in the following way: it consists of households measured in

- 2005–2009 - with rapidly decreasing income,
- 2009–2011 - with very low disposable income,
- 2011–2016 - with steeply increasing income,

but not necessarily following this trend for the whole span of 12 years, analogously for the golden cluster. Hence, these two clusters do not represent two of the typical outcome evolutions of a Czech household. This is the reason why we consider them rather an overfitting issue than actual clusters worthy of exploration. Which leaves us with $K = 4$ solution being the most suitable for the overall interpretation.

Households in the violet cluster with exceptionally high income can also always afford to pay for one week of holiday abroad and do not find the housing cost to be a really heavy burden. On the other hand, turquoise cluster represents households of completely reverse characteristics - very low disposable income, inability to pay for a week holiday abroad and almost all of them struggle with payments for housing. The other two remaining clusters (red and green) share very similar and ordinary evolution of total disposable income, but can still be distinguished. See the proportions of categorical outcomes changing in time, especially the years 2010 and 2011 when the red cluster has the lowest percentage of households able to pay for week holiday, while the green cluster

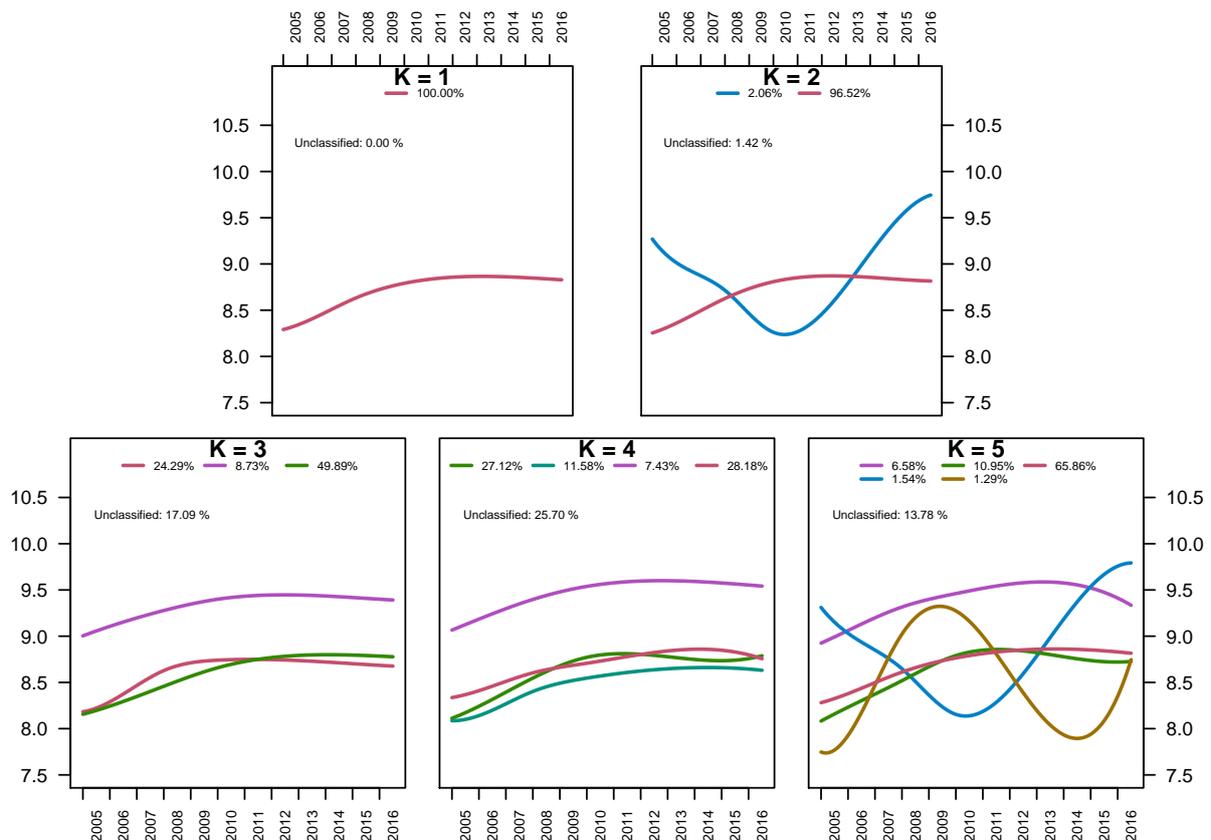


Fig. 6: Spline curves for logarithm of the total disposable household income of unit weighted family size for different choice of the number of clusters K .

has the largest. Moreover, the evolution of proportions in both categorical outcomes is reversely mirrored, when one cluster thrives the other struggles and vice versa. It almost seems like that one big 60 % cluster of average households was divided in half based on the undergoing positive or negative changes at certain periods of time. This division was allowed by our spline parametrization and the 4-year rotational panel invoked by the EU-SILC study.

7 Conclusion

In this paper we faced the problem of joint modelling of longitudinally measured numeric, ordinal and binary outcomes. We proposed to use multivariate linear mixed effects model on numeric and latent numeric outcomes corresponding to the categorical ones by exploiting the thresholding concept. Supposing all random effects to follow a joint normal distribution we enabled the outcomes to be related as it is common situation in real data analyses. On top of that, we enriched the model by construction of a mixture of such models allowing us to cluster individuals into several groups of various patterns in terms of time evolution or the variance covariance structure. By setting reasonable prior distributions for model and auxiliary parameters a hierarchical model was created and fully Bayesian approach was adopted and then executed using Gibbs sampling as one of the methods of MCMC. [Appendix](#) provides extensive derivations of full conditioned distributions used within the sampling mechanism. Sampled parameter values were not only used for the model estimation but for the calculation of classification probabilities as well, detailed derivation of which was provided in [Section 4.7](#).

Proposed model and estimation method were tested in a simulation study with the aim to examine the ability to properly estimate model parameters and to correctly capture the patterns of each individual cluster. And indeed, we verified the consistency of parameter estimates even in the case of categorical variables modelled by simple thresholding. On the other hand, certain issues with the rate of convergence of variance matrix Σ or, especially, thresholds of latent numeric counterparts of ordinal variables have appeared. This motivates us to improve or even replace current Gibbs sampling with some more advanced MCMC techniques. That may include even automated selection of the total number of clusters K , where the inspiration for the future work comes from [Neal \(2000\)](#) and his use of Dirichlet process. Variance matrix Σ of random effects needs a careful attention as well as its dimension

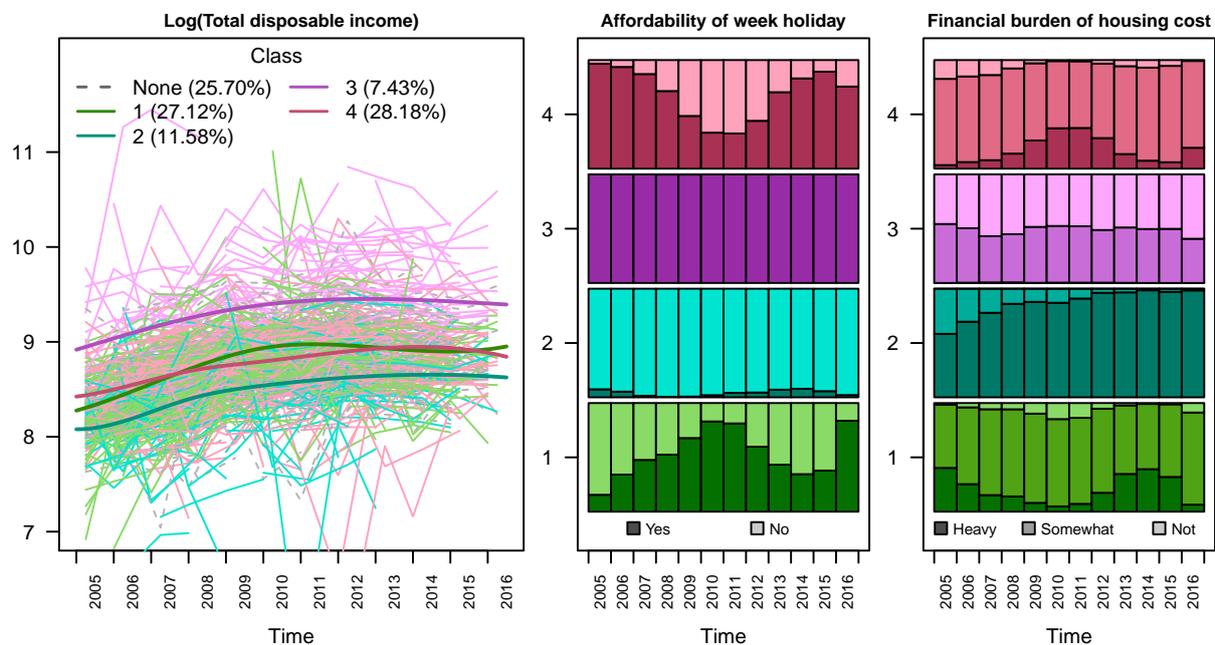


Fig. 7: Longitudinal profiles of numeric, binary and ordinal outcomes of $n = 1000$ randomly selected Czech households. Bold curves on the left represent the estimated conditional expectation of response within $K = 4$ discovered groups for a household of unit weighted family size. Categorical outcomes are presented by the proportions of categories in each year separately for the discovered groups. Some households remain unclassified.

risers with the number of modelled outcomes and the complexity of random effects structure. It may prove useful to abandon the complete generality and replace it with some commonly used block structures of variance matrices.

So far we examined the properties under fixed and minimal number of outcomes. Next work should be focused on the case of much higher number of measured outcomes and possibly on their relevance towards clustering using, for example, methods presented by [Raftery and Dean \(2006\)](#). The variable selection process could also be extended to the regression part part of the model. For example, in the EU-SILC database each household has several potential characteristics (family size, type of dwelling, number of rooms, degree of urbanization, region, country, gender, age or education of the head of a household, . . .), influence of which on outcomes and subsequent clustering may be of interest.

Regarding the real data analysis, we successfully managed to discover several different patterns in the evolution of total disposable income, affordability of a week holiday abroad and self-evaluation of difficulty to pay for housing. Minorities of extremely high (7.43 %) or low (11.58 %) life standard are easily distinguishable unlike the other mid-class households of similar income and categorical proportions but different periods of increasing and decreasing tendencies. Using more than 4-cluster solution separates out households undergoing a huge progress or recession during a certain period of time. Though these findings may be of some interest, the corresponding patterns as a whole are unrealistic due to rotational design of the study, which demands only observations from four consecutive years, and hence, these clusters are irrelevant from the realistic point of view.

The whole estimation process was implemented completely from scratch using \mathbb{R} for data preparation and processing while \mathbb{C} functions were called during the sampling phase in hope to reduce the computation time. This combination truly proved to be more efficient than pure \mathbb{R} implementation and resulted in a matter of minutes. The computation time, however, mainly depends on the number of subjects n - the larger it is, the higher number of latent parameters is to be sampled. For example, each of the cluster indicators U_i requires computation of K full conditional probabilities (38), which is still easily manageable as the latent numeric outcomes are at our disposal. This however, does not hold true for the classification probabilities (14) computed right after the sampling using an additional \mathbb{C} function within \mathbb{R} environment, since an integration over latent numeric outcomes need to be performed. The temporarily best solution (in terms of computation time and accuracy) involved calling a \mathbb{C} version of `pmvnorm` function from `mvtnorm` package which itself uses MCMC principles. Triggering this iterative process K -times for each individual and each sampled set of parameters takes a heavy toll. Therefore, some thinning methods were applied in order to obtain results in a still reasonable time.

Acknowledgement

This research was supported by the Czech Science Foundation (GAČR) grant 19-00015S.

Appendix

A Full-conditioned distributions in Gibbs sampling

In this section we derive full-conditioned distributions for all parameters $\boldsymbol{\psi} \in \boldsymbol{\Psi}$ one by one. All derivations are based on viewing (26) as a function of parameter $\boldsymbol{\psi}$. The function on the right hand side of (26) can be decomposed into the following products:

$$\begin{aligned}
p(\boldsymbol{\psi} | \mathbb{Y}, \mathcal{C}; \boldsymbol{\Psi}_{-\boldsymbol{\psi}}, \boldsymbol{\zeta}) &\propto \prod_{i=1}^n p\left(\mathbb{Y}_i^{\text{OB}} | \mathbb{Y}_i^{\star, \text{OB}}; \boldsymbol{\gamma}\right) \cdot \prod_{i=1}^n p\left(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{\star, \text{OB}} | \mathcal{C}_i, \mathbf{b}_i; \boldsymbol{\beta}^{(U_i)}, \boldsymbol{\tau}^{(U_i)}\right) \cdot \\
&\quad \cdot \prod_{i=1}^n p\left(\mathbf{b}_i | \boldsymbol{\mu}^{(U_i)}, \boldsymbol{\Sigma}^{-(U_i)}\right) \cdot \prod_{i=1}^n p(U_i | \mathbf{w}) \cdot \\
&\quad \cdot p(\mathbf{w} | \boldsymbol{\alpha}) \cdot p\left(\boldsymbol{\gamma} | \boldsymbol{\gamma}_1^{\text{O}}, r^{\text{O}} \in \mathcal{R}^{\text{O}}\right) \cdot p(\boldsymbol{\beta} | \boldsymbol{\tau}; \boldsymbol{\beta}_0, \mathbb{D}) \cdot p(\boldsymbol{\tau} | a_1, a_2) \cdot \\
&\quad \cdot p(\boldsymbol{\mu} | \boldsymbol{\tau}_{\text{R}}; \boldsymbol{\mu}_0) \cdot p(\boldsymbol{\tau}_{\text{R}} | a_3, a_4) \cdot p(\boldsymbol{\Sigma}^{-1} | \mathbb{Q}^{-1}; \mathbf{v}_0) \cdot p\left(\mathbb{Q}^{-1} | \mathbb{D}^{\mathbb{Q}}, \mathbf{v}_1\right). \quad (36)
\end{aligned}$$

Derivations are made under the assumption that parameters $\boldsymbol{\beta}$, $\boldsymbol{\tau}$, $\boldsymbol{\mu}$, $\boldsymbol{\tau}_{\text{R}}$, $\boldsymbol{\Sigma}^{-1}$ and \mathbb{Q}^{-1} are all group-specific. Similar derivations (with corresponding changes) can be made in the case of chosen subset of group-specific parameters. Note that if $\boldsymbol{\tau}_{\text{R}}$ and \mathbb{Q}^{-1} are group-specific, then $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}^{-1}$ must be group-specific as well.

A.1 Probabilities \mathbf{w}

Parameter \mathbf{w} of prior probabilities of being classified into certain categories, i.e. $w_k = \mathbb{P}(U_i = k)$, appear only in $p(U_i|\mathbf{w})$ and its prior distribution $p(\mathbf{w}|\boldsymbol{\alpha})$, therefore:

$$p(\mathbf{w}|\mathbb{Y}, \mathcal{C}; \boldsymbol{\Psi}_{-\mathbf{w}}; \boldsymbol{\zeta}) \propto \prod_{i=1}^n p(U_i|\mathbf{w}) \cdot p(\mathbf{w}|\boldsymbol{\alpha})$$

$$p(\mathbf{w}|\mathbf{U}; \boldsymbol{\alpha}) \propto \prod_{k=1}^K w_k^{\sum_{i=1}^n \mathbb{1}(U_i=k)} \cdot \prod_{k=1}^K w_k^{\alpha_k-1} = \prod_{k=1}^K w_k^{n^k(\mathbf{U})+\alpha_k-1},$$

where $n^k(\mathbf{U})$ is the total number of appearances of value k among all current group-allocation indicators $\mathbf{U} = \{U_i, i = 1, \dots, n\}$, i.e. the total number of subjects (from n possible) currently belonging to group k . Comparing this pdf with (19) we recognize the shape of pdf of Dirichlet distribution, thus,

$$\mathbf{w}|\mathbb{Y}, \mathcal{C}; \boldsymbol{\Psi}_{-\mathbf{w}}; \boldsymbol{\zeta} \sim \text{Dir}_K(\mathbf{n}(\mathbf{U}) + \boldsymbol{\alpha}), \quad (37)$$

where $\mathbf{n}(\mathbf{U}) = (n^1(\mathbf{U}), \dots, n^K(\mathbf{U}))^\top$.

A.2 Group-allocation indicators U_i

According to (36) the group-allocation indicator U_i appears only in its prior distribution $U_i|\mathbf{w}$ and at places, where it selects the corresponding group-specific parameter:

$$p(U_i|\mathbb{Y}, \mathcal{C}; \boldsymbol{\Psi}_{-U_i}; \boldsymbol{\zeta}) \propto p\left(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{\text{*,OB}} \mid \mathcal{C}_i, \mathbf{b}_i; \boldsymbol{\beta}^{(U_i)}, \boldsymbol{\tau}^{(U_i)}\right) \cdot p\left(\mathbf{b}_i \mid \boldsymbol{\mu}^{(U_i)}, \boldsymbol{\Sigma}^{-(U_i)}\right) \cdot p(U_i|\mathbf{w}).$$

U_i only attains values $k \in \{1, \dots, K\}$, therefore, we aim to calculate full-conditioned probability that i th subject is allocated in group k :

$$\mathbb{P}\left(U_i = k \mid \mathbb{Y}_i^{\text{N}}, \mathcal{C}_i; \mathbb{Y}_i^{\text{*,OB}}, \mathbf{b}_i; \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}, \mathbf{w}\right) \propto w_k \cdot \prod_{r^{\text{N}} \in \mathcal{R}^{\text{Num}}} \left(\tau_{r^{\text{N}}}^{(k)}\right)^{\frac{n_i}{2}} \cdot \left|\boldsymbol{\Sigma}^{-(k)}\right|^{\frac{1}{2}} \cdot \exp\left\{-\frac{1}{2} \sum_{r^{\text{N}} \in \mathcal{R}^{\text{Num}}} \sum_{j=1}^{n_i} \tau_{r^{\text{N}}}^{(k)} \left(y_{i,j}^{r^{\text{N}}} - \eta_{i,j}^{(k),r^{\text{N}}}\right)^2 - \frac{1}{2} \sum_{r^{\text{OB}} \in \mathcal{R}^{\text{OB}}} \sum_{j=1}^{n_i} \left(y_{i,j}^{*,r^{\text{OB}}} - \eta_{i,j}^{(k),r^{\text{OB}}}\right)^2\right\} \cdot \exp\left\{-\frac{1}{2} \left(\mathbf{b}_i - \boldsymbol{\mu}^{(k)}\right)^\top \boldsymbol{\Sigma}^{-(k)} \left(\mathbf{b}_i - \boldsymbol{\mu}^{(k)}\right)\right\}, \quad (38)$$

where $\eta_{i,j}^{(k),r} = \left(\mathbf{x}_{i,j}^r\right)^\top \boldsymbol{\beta}^{(k),r} + \left(\mathbf{z}_{i,j}^r\right)^\top \mathbf{b}_i$ is the linear predictor of j -th observation of outcome $r \in \mathcal{R}$ of i th subject in group k .

A.3 Latent numeric variables $\mathbb{Y}_i^{\text{*,OB}}$

Latent numeric outcomes $\mathbb{Y}_i^{\text{*,OB}}$ for actually measured ordinal and binary outcomes \mathbb{Y}_i^{OB} appear only in the thresholding procedure and multivariate LME for both \mathbb{Y}_i^{N} and $\mathbb{Y}_i^{\text{*,OB}}$:

$$p\left(\mathbb{Y}_i^{\text{*,OB}} \mid \mathbb{Y}, \mathcal{C}; \boldsymbol{\Psi}_{-\mathbb{Y}_i^{\text{*,OB}}}; \boldsymbol{\zeta}\right) \propto p\left(\mathbb{Y}_i^{\text{OB}} \mid \mathbb{Y}_i^{\text{*,OB}}; \boldsymbol{\gamma}\right) \cdot p\left(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{\text{*,OB}} \mid \mathcal{C}_i, \mathbf{b}_i; \boldsymbol{\beta}^{(U_i)}, \boldsymbol{\tau}^{(U_i)}\right).$$

From (7) we see that for all $r^{\text{OB}} \in \mathcal{R}^{\text{OB}}$ and $j = 1, \dots, n_i$ are $Y_{i,j}^{*,r^{\text{OB}}}$ independently distributed. Ignoring the thresholding concept $Y_{i,j}^{*,r^{\text{OB}}}$ would follow $\text{N}\left(\eta_{i,j}^{(U_i),r^{\text{OB}}}, 1\right)$, however, corresponding density is now limited by indicator $\mathbb{1}\left(\mathcal{Y}_l^{\text{OB}}, \mathcal{Y}_{l+1}^{\text{OB}}\right]\left(y_i^{*,\text{OB}}\right)$, where $l = y_{i,j}^{\text{OB}}$. Therefore, the full-conditioned distribution is truncated normal distribution on interval $\left(\mathcal{Y}_l^{\text{OB}}, \mathcal{Y}_{l+1}^{\text{OB}}\right]$:

$$Y_{i,j}^{*,r^{\text{OB}}} \mid Y_{i,j}^{\text{OB}} = l, \boldsymbol{\gamma} \sim \text{TN}\left(\eta_{i,j}^{(U_i),r^{\text{OB}}}, 1, \mathcal{Y}_l^{\text{OB}}, \mathcal{Y}_{l+1}^{\text{OB}}\right). \quad (39)$$

A.4 Thresholds $\boldsymbol{\gamma}$

Parameter $\boldsymbol{\gamma}$ affects (36) only in the thresholding phase and in prior distribution of $\boldsymbol{\gamma}$:

$$p\left(\boldsymbol{\gamma} \mid \mathbb{Y}, \mathcal{C}; \boldsymbol{\Psi}_{-\boldsymbol{r}^{\text{O}}}; \boldsymbol{\zeta}\right) \propto \prod_{i=1}^n p\left(\mathbb{Y}_i^{\text{OB}} \mid \mathbb{Y}_i^{\text{OB}}; \boldsymbol{\gamma}\right) \cdot p\left(\boldsymbol{\gamma} \mid \boldsymbol{\gamma}_1^{\text{O}}, \boldsymbol{r}^{\text{O}} \in \mathcal{R}^{\text{O}}\right)$$

Let us consider ordinal outcome $\boldsymbol{r}^{\text{O}} \in \mathcal{R}^{\text{Ord}}$ and corresponding set of thresholds: $-\infty = \boldsymbol{\gamma}_0^{\text{O}}, \boldsymbol{\gamma}_1^{\text{O}}, \boldsymbol{\gamma}^{\text{O}}, \boldsymbol{\gamma}_{L^{\text{O}}}^{\text{O}} = \infty$. Let $\mathcal{Y}_l^{\boldsymbol{r}^{\text{O}}}$ be the set of all latent numeric outcomes $Y_{i,j}^{*,\boldsymbol{r}^{\text{O}}}$ such that the truly measured ordinal category is $l = 0, \dots, L^{\text{O}} - 1$, i.e.

$$\mathcal{Y}_l^{\boldsymbol{r}^{\text{O}}} = \left\{ Y_{i,j}^{*,\boldsymbol{r}^{\text{O}}} : Y_{i,j}^{\text{O}} = l, \quad i = 1, \dots, n, \quad j = 1, \dots, n_i \right\},$$

which is assumed to be non-empty (all levels of outcome L^{O} are attained at least once). The latent numeric variables had to be generated according to the thresholding concept, therefore, the following inequalities hold:

$$-\infty < \underset{\in \mathcal{Y}_0^{\boldsymbol{r}^{\text{O}}}}{y_0} < \underset{\in \mathcal{Y}_1^{\boldsymbol{r}^{\text{O}}}}{\boldsymbol{\gamma}_1^{\text{O}}} < \underset{\in \mathcal{Y}_1^{\boldsymbol{r}^{\text{O}}}}{y_1} < \underset{\in \mathcal{Y}_2^{\boldsymbol{r}^{\text{O}}}}{\boldsymbol{\gamma}_2^{\text{O}}} < \underset{\in \mathcal{Y}_2^{\boldsymbol{r}^{\text{O}}}}{y_2} < \dots < \underset{\in \mathcal{Y}_{L^{\text{O}}-1}^{\boldsymbol{r}^{\text{O}}}}{\boldsymbol{\gamma}_{L^{\text{O}}-1}^{\text{O}}} < \underset{\in \mathcal{Y}_{L^{\text{O}}-1}^{\boldsymbol{r}^{\text{O}}}}{y_{L^{\text{O}}-1}} < \infty.$$

Thus, under the uniform prior (20) for $\boldsymbol{\gamma}^{\text{O}}$ we get that individual thresholds $\boldsymbol{\gamma}_l^{\text{O}}$ are uniformly distributed on intervals given by maxima and minima of corresponding sets:

$$\boldsymbol{\gamma}_l^{\text{O}} \mid \boldsymbol{Y}^{\boldsymbol{r}^{\text{O}}}; \boldsymbol{Y}^{*,\boldsymbol{r}^{\text{O}}} \sim \text{Unif} \left[\max_{y \in \mathcal{Y}_{l-1}^{\boldsymbol{r}^{\text{O}}}} y, \min_{y \in \mathcal{Y}_l^{\boldsymbol{r}^{\text{O}}}} y \right], \quad l = 1, \dots, L^{\text{O}}. \quad (40)$$

A.5 Precision parameters $\boldsymbol{\tau}$

Parameters $\boldsymbol{\tau} = \left\{ \boldsymbol{\tau}_{r^{\text{N}}}^{(k)} : k = 1, \dots, K, \boldsymbol{r}^{\text{N}} \in \mathcal{R}^{\text{Num}} \right\}$ are the inverse variance of errors of supposed LME models over numeric outcomes. The right-hand side of (36) includes $\boldsymbol{\tau}$ only in three factors - supposed LME for \mathbb{Y}_i^{N} and prior distribution of $(\boldsymbol{\beta}, \boldsymbol{\tau})$:

$$p\left(\boldsymbol{\tau} \mid \mathbb{Y}, \mathcal{C}; \boldsymbol{\Psi}_{-\boldsymbol{\tau}}; \boldsymbol{\zeta}\right) \propto \prod_{i=1}^n p\left(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{\text{OB}} \mid \mathcal{C}_i, \boldsymbol{b}_i; \boldsymbol{\beta}^{(U_i)}, \boldsymbol{\tau}^{(U_i)}\right) \cdot p\left(\boldsymbol{\beta} \mid \boldsymbol{\tau}; \boldsymbol{\beta}_0, \mathbb{D}\right) \cdot p\left(\boldsymbol{\tau} \mid a_1, a_2\right)$$

From the structure of (7), (22) and (23) we see that individual $\boldsymbol{\tau}_{r^{\text{N}}}^{(k)}$ are distributed independently of each other (given all other information and parameters):

$$p\left(\boldsymbol{\tau}_{r^{\text{N}}}^{(k)} \mid \boldsymbol{Y}^{r^{\text{N}}}, \mathcal{C}^{r^{\text{N}}}; \boldsymbol{U}, \boldsymbol{b}^{r^{\text{N}}}; \boldsymbol{\beta}^{(k),r^{\text{N}}}; \boldsymbol{\beta}_0^{r^{\text{N}}}, \mathbb{D}^{r^{\text{N}}}, a_1, a_2\right) \propto \left(\boldsymbol{\tau}_{r^{\text{N}}}^{(k)}\right)^{\frac{1}{2} \sum_{i=1}^n n_i \mathbb{1}(U_i=k) + \frac{1}{2} d_{r^{\text{N}}}^{\text{F}} + a_1 - 1} \cdot \exp \left\{ -\boldsymbol{\tau}_{r^{\text{N}}}^{(k)} \left[\frac{1}{2} \sum_{i \in \mathcal{N}_k(\boldsymbol{U})} \sum_{j=1}^{n_i} \left(y_{i,j}^{r^{\text{N}}} - \eta_{i,j}^{(k),r^{\text{N}}} \right)^2 + \frac{1}{2} \sum_{j=1}^{d_{r^{\text{N}}}^{\text{F}}} \frac{\left(\boldsymbol{\beta}_j^{(k),r} - \boldsymbol{\beta}_{0,j}^r \right)^2}{d_{j,j}^r} + a_2 \right] \right\}, \quad (41)$$

where $\mathcal{N}_k(\boldsymbol{U}) = \{i : U_i = k, i = 1, \dots, n\}$ is a set of subjects currently belonging to group k . For $\boldsymbol{Y}^{r^{\text{N}}}, \mathcal{C}^{r^{\text{N}}}$ and current values of $\boldsymbol{U}, \boldsymbol{b}^{r^{\text{N}}}$ and $\boldsymbol{\beta}^{(k)}$ let us denote

$$\begin{aligned} \tilde{a}_1^{(k),r^{\text{N}}} &= \frac{1}{2} \sum_{i \in \mathcal{N}_k(\boldsymbol{U})} n_i + \frac{d_{r^{\text{N}}}^{\text{F}}}{2} + a_1, \\ \tilde{a}_2^{(k),r^{\text{N}}} &= \frac{1}{2} \sum_{i \in \mathcal{N}_k(\boldsymbol{U})} \sum_{j=1}^{n_i} \left(y_{i,j}^{r^{\text{N}}} - \eta_{i,j}^{(k),r^{\text{N}}} \right)^2 + \frac{1}{2} \sum_{j=1}^{d_{r^{\text{N}}}^{\text{F}}} \frac{\left(\boldsymbol{\beta}_j^{(k),r} - \boldsymbol{\beta}_{0,j}^r \right)^2}{d_{j,j}^r} + a_2. \end{aligned}$$

Comparing (41) with (23) we see that

$$\boldsymbol{\tau}_{r^{\text{N}}}^{(k)} \mid \boldsymbol{Y}^{r^{\text{N}}}, \mathcal{C}^{r^{\text{N}}}; \boldsymbol{U}, \boldsymbol{b}^{r^{\text{N}}}; \boldsymbol{\beta}^{(k),r^{\text{N}}}; \boldsymbol{\beta}_0^{r^{\text{N}}}, \mathbb{D}^{r^{\text{N}}}, a_1, a_2 \sim \Gamma\left(\tilde{a}_1^{(k),r^{\text{N}}}, \tilde{a}_2^{(k),r^{\text{N}}}\right) \quad (42)$$

independently for each $r^{\text{N}} \in \mathcal{R}^{\text{Num}}$ and $k = 1, \dots, K$.

A.6 Fixed effects $\boldsymbol{\beta}$

Fixed effects $\boldsymbol{\beta}$ appear only in LME model specification and prior distribution:

$$p(\boldsymbol{\beta} | \mathbb{Y}, \mathcal{C}; \boldsymbol{\Psi}, \boldsymbol{\beta}; \boldsymbol{\zeta}) \propto \prod_{i=1}^n p\left(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{\text{*,OB}} \mid \mathcal{C}_i, \mathbf{b}_i; \boldsymbol{\beta}^{(U_i)}, \boldsymbol{\tau}^{(U_i)}\right) \cdot p(\boldsymbol{\beta} | \boldsymbol{\tau}; \boldsymbol{\beta}_0, \mathbb{D}),$$

which can be decomposed for individual outcomes $r \in \mathcal{R}$ and $k = 1, \dots, K$ as follows:

$$p\left(\boldsymbol{\beta}^{(k),r} \mid \mathbf{Y}^r, \mathcal{C}^r; \mathbf{U}, \mathbf{b}^r; \boldsymbol{\tau}_r^{(k)}; \boldsymbol{\beta}_0, \mathbb{D}^r\right) \propto \exp\left\{-\frac{\boldsymbol{\tau}_r^{(k)}}{2} \left(\boldsymbol{\beta}^{(k),r} - \boldsymbol{\beta}_0^r\right)^\top [\mathbb{D}^r]^{-1} \left(\boldsymbol{\beta}^{(k),r} - \boldsymbol{\beta}_0^r\right)\right\} \\ \cdot \exp\left\{-\frac{\boldsymbol{\tau}_r^{(k)}}{2} \left(\tilde{\mathbf{y}}_{\mathcal{N}_k^r(\mathbf{U})}^r - \mathbb{X}_{\mathcal{N}_k^r(\mathbf{U})}^r \boldsymbol{\beta}^{(k),r}\right)^\top \left(\tilde{\mathbf{y}}_{\mathcal{N}_k^r(\mathbf{U})}^r - \mathbb{X}_{\mathcal{N}_k^r(\mathbf{U})}^r \boldsymbol{\beta}^{(k),r}\right)\right\},$$

where notation $\bullet_{\mathcal{N}_k^r(\mathbf{U})}$ restricts given parameter \bullet to the subset of subjects in group k :

$$\mathbb{X}_{\mathcal{N}_k^r(\mathbf{U})}^r = \begin{pmatrix} \vdots \\ \mathbb{X}_i^r \\ \vdots \end{pmatrix}, i \in \mathcal{N}_k^r(\mathbf{U}), \quad \tilde{\mathbf{y}}_{\mathcal{N}_k^r(\mathbf{U})}^r = \begin{cases} \left[(\mathbf{y}_i^r - \mathbb{Z}_i^r \mathbf{b}_i^r)^\top, i \in \mathcal{N}_k^r(\mathbf{U})\right]^\top, & \text{if } r \in \mathcal{R}^{\text{Num}}, \\ \left[(\mathbf{y}_i^{*,r} - \mathbb{Z}_i^r \mathbf{b}_i^r)^\top, i \in \mathcal{N}_k^r(\mathbf{U})\right]^\top, & \text{if } r \in \mathcal{R}^{\text{OB}}. \end{cases}$$

Using basic algebraic operations and ignoring several multiplicative constants we can rewrite probability density function of full-conditioned distribution of $\boldsymbol{\beta}^{(k),r}$ into:

$$p\left(\boldsymbol{\beta}^{(k),r} \mid \mathbf{Y}^r, \mathcal{C}^r; \mathbf{U}, \mathbf{b}^r; \boldsymbol{\tau}_r^{(k)}; \boldsymbol{\beta}_0, \mathbb{D}^r\right) \propto \\ \exp\left\{-\frac{\boldsymbol{\tau}_r^{(k)}}{2} \left(\boldsymbol{\beta}^{(k),r} - \tilde{\boldsymbol{\beta}}^{(k),r}\right)^\top \left[\left(\mathbb{X}_{\mathcal{N}_k^r(\mathbf{U})}^r\right)^\top \mathbb{X}_{\mathcal{N}_k^r(\mathbf{U})}^r + (\mathbb{D}^r)^{-1}\right] \left(\boldsymbol{\beta}^{(k),r} - \tilde{\boldsymbol{\beta}}^{(k),r}\right)\right\}, \quad (43)$$

where

$$\tilde{\boldsymbol{\beta}}^{(k),r} = \left[\left(\mathbb{X}_{\mathcal{N}_k^r(\mathbf{U})}^r\right)^\top \mathbb{X}_{\mathcal{N}_k^r(\mathbf{U})}^r + (\mathbb{D}^r)^{-1}\right]^{-1} \left(\left(\mathbb{X}_{\mathcal{N}_k^r(\mathbf{U})}^r\right)^\top \tilde{\mathbf{y}}_{\mathcal{N}_k^r(\mathbf{U})}^r + (\mathbb{D}^r)^{-1} \boldsymbol{\beta}_0^r\right),$$

which compared to pdf of multivariate normal distribution yields

$$\boldsymbol{\beta}^{(k),r} \mid \mathbf{Y}^r, \mathcal{C}^r; \mathbf{U}, \mathbf{b}^r; \boldsymbol{\tau}_r^{(k)}; \boldsymbol{\beta}_0, \mathbb{D}^r \sim \text{N}_{d_{\mathbb{F}}} \left(\tilde{\boldsymbol{\beta}}^{(k),r}, \frac{\left[\left(\mathbb{X}_{\mathcal{N}_k^r(\mathbf{U})}^r\right)^\top \mathbb{X}_{\mathcal{N}_k^r(\mathbf{U})}^r + (\mathbb{D}^r)^{-1}\right]^{-1}}{\boldsymbol{\tau}_r^{(k)}} \right). \quad (44)$$

A.7 Prior precisions $\boldsymbol{\tau}_{\mathbf{R}}$ for $\boldsymbol{\mu}$

Parameter $\boldsymbol{\tau}_{\mathbf{R}}$ serves as an auxiliary parameter for specifying prior distribution of $\boldsymbol{\mu}$, see section 4.6.1. The derivation of full-conditioned distribution of this parameter is solely based on combining probability distribution functions in (24). Therefore,

$$p(\boldsymbol{\tau}_{\mathbf{R}} \mid \boldsymbol{\mu}; \boldsymbol{\mu}_0, a_3, a_4) \propto \prod_{k=1}^K \prod_{j=1}^{d^{\mathbf{R}}} \left(\boldsymbol{\tau}_{\mathbf{R},j}^{(k)}\right)^{a_3 + \frac{1}{2} - 1} \exp\left\{-\boldsymbol{\tau}_{\mathbf{R},j}^{(k)} \left[a_4 + \frac{1}{2} \left(\boldsymbol{\mu}_j^{(k)} - \boldsymbol{\mu}_{0,j}^{(k)}\right)^2\right]\right\} \quad (45)$$

and

$$\boldsymbol{\tau}_{\mathbf{R},j}^{(k)} \mid \boldsymbol{\mu}_j^{(k)}; \boldsymbol{\mu}_{0,j}^{(k)}, a_3, a_4 \sim \Gamma\left(a_3 + \frac{1}{2}, a_4 + \frac{1}{2} \left(\boldsymbol{\mu}_j^{(k)} - \boldsymbol{\mu}_{0,j}^{(k)}\right)^2\right) \quad (46)$$

independently for all $j = 1, \dots, d^{\mathbf{R}}$ and $k = 1, \dots, D$.

A.8 Prior expected values $\boldsymbol{\mu}$ for \mathbf{b}

Parameter $\boldsymbol{\mu}$ consists of all possible expected values $\boldsymbol{\mu}^{(k)}$ of random effects \mathbf{b}_i in all groups $k = 1, \dots, K$. The right-hand side of (36) is in the case of this parameter simplified into

$$p(\boldsymbol{\mu} | \mathbb{Y}, \mathcal{C}; \boldsymbol{\Psi}_{-\boldsymbol{\mu}}; \boldsymbol{\zeta}) \propto \prod_{i=1}^n p(\mathbf{b}_i | \boldsymbol{\mu}^{(U_i)}, \boldsymbol{\Sigma}^{-(U_i)}) \cdot p(\boldsymbol{\mu} | \boldsymbol{\tau}_R; \boldsymbol{\mu}_0).$$

From the product across all subjects for given group $k = 1, \dots, K$ we extract only those factors that correspond to subjects within k -th group, i.e. $\mathcal{N}_k(\mathbf{U})$. By performing several algebraic operations and ignoring multiplicative constants, we obtain

$$\begin{aligned} p(\boldsymbol{\mu}^{(k)} | \mathbf{U}, \mathbf{b}; \boldsymbol{\Sigma}^{-(k)}; \boldsymbol{\tau}_R^{(k)}; \boldsymbol{\mu}_0^{(k)}) &\propto \prod_{i \in \mathcal{N}_k(\mathbf{U})} p(\mathbf{b}_i | \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{-(k)}) \cdot p(\boldsymbol{\mu}^{(k)} | \boldsymbol{\tau}_R^{(k)}; \boldsymbol{\mu}_0^{(k)}) \\ &\propto \exp \left\{ -\frac{1}{2} \sum_{i \in \mathcal{N}_k(\mathbf{U})} (\mathbf{b}_i - \boldsymbol{\mu}^{(k)})^\top \boldsymbol{\Sigma}^{-(k)} (\mathbf{b}_i - \boldsymbol{\mu}^{(k)}) - \frac{1}{2} (\boldsymbol{\mu}^{(k)} - \boldsymbol{\mu}_0^{(k)})^\top \text{diag}(\boldsymbol{\tau}_R^{(k)}) (\boldsymbol{\mu}^{(k)} - \boldsymbol{\mu}_0^{(k)}) \right\} \quad (47) \\ &\propto \exp \left\{ -\frac{1}{2} (\boldsymbol{\mu}^{(k)} - \tilde{\boldsymbol{\mu}}^{(k)})^\top \left[n^k(\mathbf{U}) \boldsymbol{\Sigma}^{-(k)} + \text{diag}(\boldsymbol{\tau}_R^{(k)}) \right] (\boldsymbol{\mu}^{(k)} - \tilde{\boldsymbol{\mu}}^{(k)}) \right\}, \end{aligned}$$

where

$$\tilde{\boldsymbol{\mu}}^{(k)} = \left[n^k(\mathbf{U}) \boldsymbol{\Sigma}^{-(k)} + \text{diag}(\boldsymbol{\tau}_R^{(k)}) \right]^{-1} \left(n^k(\mathbf{U}) \boldsymbol{\Sigma}^{-(k)} \underbrace{\frac{1}{n^k(\mathbf{U})} \sum_{i \in \mathcal{N}_k(\mathbf{U})} \mathbf{b}_i}_{\tilde{\mathbf{b}}^k(\mathbf{U})} + \text{diag}(\boldsymbol{\tau}_R^{(k)}) \boldsymbol{\mu}_0 \right)$$

leading to the following full-conditioned distribution

$$\boldsymbol{\mu}^{(k)} | \mathbf{U}, \mathbf{b}; \boldsymbol{\Sigma}^{-(k)}; \boldsymbol{\tau}_R^{(k)}; \boldsymbol{\mu}_0^{(k)} \sim N_{d^R} \left(\tilde{\boldsymbol{\mu}}^{(k)}, \left[n^k(\mathbf{U}) \boldsymbol{\Sigma}^{-(k)} + \text{diag}(\boldsymbol{\tau}_R^{(k)}) \right]^{-1} \right) \quad (48)$$

independently for all $k = 1, \dots, K$.

A.9 Prior scale matrices \mathbb{Q}^{-1} for $\boldsymbol{\Sigma}^{-1}$

Parameter \mathbb{Q}^{-1} is the set of auxiliary parameters that makes prior distribution of $\boldsymbol{\Sigma}^{-1}$ more flexible within Gibbs sampler. The right-hand side of (36) shrinks into

$$p(\mathbb{Q}^{-1} | \mathbb{Y}, \mathcal{C}; \boldsymbol{\Psi}_{-\mathbb{Q}^{-1}}; \boldsymbol{\zeta}) \propto p(\boldsymbol{\Sigma}^{-1} | \mathbb{Q}^{-1}; \mathbf{v}_0) \cdot p(\mathbb{Q}^{-1} | \mathbb{D}^{\mathbb{Q}}, \mathbf{v}_1).$$

Combining the two probability density functions in (25) we get

$$p(\mathbb{Q}^{-(k)} | \boldsymbol{\Sigma}^{-(k)}; \mathbf{v}_0, \mathbf{v}_1, \mathbb{D}^{\mathbb{Q}}) \propto |\mathbb{Q}^{-(k)}|^{-\frac{\mathbf{v}_0 + \mathbf{v}_1 - d^R - 1}{2}} \exp \left\{ -\text{Tr} \left[\left(\boldsymbol{\Sigma}^{-(k)} + (\mathbb{D}^{\mathbb{Q}})^{-1} \right) \mathbb{Q}^{-(k)} \right] \right\}, \quad (49)$$

which compared to (25) resembles pdf of Wishart distribution. Therefore,

$$\mathbb{Q}^{-(k)} | \boldsymbol{\Sigma}^{-(k)}; \mathbf{v}_0, \mathbf{v}_1, \mathbb{D}^{\mathbb{Q}} \sim W_{d^R} \left(\left[\boldsymbol{\Sigma}^{-(k)} + (\mathbb{D}^{\mathbb{Q}})^{-1} \right]^{-1}, \mathbf{v}_0 + \mathbf{v}_1 \right) \quad (50)$$

independently for all $k = 1, \dots, K$.

A.10 Prior inverse covariance matrices Σ^{-1} for random effects \mathbf{b}

Parameter Σ^{-1} is the set of inverse covariance matrices for random effects \mathbf{b}_i that contributes to the right-hand side of 36 only in pdf for random effects and in prior distribution of Σ^{-1} :

$$p(\Sigma^{-1} | \mathbb{Y}, \mathcal{C}; \Psi_{-\Sigma^{-1}}; \zeta) \propto \prod_{i=1}^n p(\mathbf{b}_i | \boldsymbol{\mu}^{(U_i)}, \Sigma^{-(U_i)}) \cdot p(\Sigma^{-1} | \mathbb{Q}^{-1}; \mathbf{v}_0)$$

Again, we need to separate subjects into groups $\mathcal{N}_k(\mathbf{U}), k = 1, \dots, K$ according to their current allocation indicators \mathbf{U} . Similarly as before, the equation above decomposes into K independent parts - one for each group $k = 1, \dots, K$. Considering group k the right-hand side of equation above reduces into

$$\begin{aligned} p(\Sigma^{-(k)} | \mathbf{U}, \mathbf{b}; \boldsymbol{\mu}^{(k)}, \mathbb{Q}^{-(k)}; \mathbf{v}_0) &\propto |\Sigma^{-(k)}|^{\frac{n^k(\mathbf{U}) + \mathbf{v}_0 - d_{\mathbb{R}} - 1}{2}} \\ &\cdot \exp \left\{ -\frac{1}{2} \sum_{i \in \mathcal{N}_k(\mathbf{U})} (\mathbf{b}_i - \boldsymbol{\mu}^{(k)})^\top \Sigma^{-(k)} (\mathbf{b}_i - \boldsymbol{\mu}^{(k)}) - \text{Tr} [\mathbb{Q}^{-(k)} \Sigma^{-(k)}] \right\} \\ &\propto |\Sigma^{-(k)}|^{\frac{n^k(\mathbf{U}) + \mathbf{v}_0 - d_{\mathbb{R}} - 1}{2}} \exp \left\{ -\text{Tr} \left[\left(\mathbb{Q}^{-(k)} + \frac{1}{2} \sum_{i \in \mathcal{N}_k(\mathbf{U})} (\mathbf{b}_i - \boldsymbol{\mu}^{(k)}) (\mathbf{b}_i - \boldsymbol{\mu}^{(k)})^\top \right) \Sigma^{-(k)} \right] \right\}, \end{aligned} \quad (51)$$

which compared to (25) resembles pdf of Wishart distribution. Therefore, independently for all $k = 1, \dots, K$

$$\Sigma^{-(k)} | \mathbf{U}, \mathbf{b}; \boldsymbol{\mu}^{(k)}, \mathbb{Q}^{-(k)}; \mathbf{v}_0 \sim \text{W}_{d_{\mathbb{R}}}(\tilde{\mathbb{Q}}^{(k)}, n^k(\mathbf{U}) + \mathbf{v}_0), \quad (52)$$

where

$$\tilde{\mathbb{Q}}^{(k)} = (\mathbb{Q}^{-(k)})^{-1} \quad \text{and} \quad \tilde{\mathbb{Q}}^{-(k)} = \mathbb{Q}^{-(k)} + \frac{1}{2} \sum_{i \in \mathcal{N}_k(\mathbf{U})} (\mathbf{b}_i - \boldsymbol{\mu}^{(k)}) (\mathbf{b}_i - \boldsymbol{\mu}^{(k)})^\top.$$

A.11 Random effects \mathbf{b}

The key role of our model is played by random effects $\mathbf{b}_i, i = 1, \dots, n$ that create linear predictors $\boldsymbol{\eta}_i^{(k),r}, k = 1, \dots, K$ and $r \in \mathcal{R}$. Probability density function of corresponding full-conditioned distribution is based on only two parts of the right-hand side of (36):

$$p(\mathbf{b} | \mathbb{Y}, \mathcal{C}; \Psi_{-\mathbf{b}}; \zeta) \propto \prod_{i=1}^n p(\mathbb{Y}_i^{\mathbb{N}}, \mathbb{Y}_i^{*,\text{OB}} | \mathcal{C}_i, \mathbf{b}_i; \boldsymbol{\beta}^{(U_i)}, \boldsymbol{\tau}^{(U_i)}) \cdot \prod_{i=1}^n p(\mathbf{b}_i | \boldsymbol{\mu}^{(U_i)}, \Sigma^{-(U_i)}).$$

Clearly, random effects \mathbf{b}_i will be distributed independently even in the full-conditioned distribution. Let us select subject i (say from group $U_i = k$) in which case its corresponding probability distribution function is of the shape

$$\begin{aligned} p(\mathbf{b}_i | \mathbb{Y}_i^{\mathbb{N}}, \mathcal{C}_i; \mathbb{Y}_i^{*,\text{OB}}, U_i; \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\mu}, \Sigma^{-1}) &\propto \prod_{r^{\mathbb{N}} \in \mathcal{R}^{\text{Num}}} \exp \left\{ -\frac{\tau_{r^{\mathbb{N}}}^{(k)}}{2} (\tilde{\mathbf{y}}_i^{r^{\mathbb{N}}} - \mathbb{Z}_i^{r^{\mathbb{N}}} \mathbf{b}_i^{r^{\mathbb{N}}})^\top (\tilde{\mathbf{y}}_i^{r^{\mathbb{N}}} - \mathbb{Z}_i^{r^{\mathbb{N}}} \mathbf{b}_i^{r^{\mathbb{N}}}) \right\} \\ &\cdot \prod_{r^{\text{OB}} \in \mathcal{R}^{\text{OB}}} \exp \left\{ -\frac{1}{2} (\tilde{\mathbf{y}}_i^{*,r^{\text{OB}}} - \mathbb{Z}_i^{r^{\text{OB}}} \mathbf{b}_i^{r^{\text{OB}}})^\top (\tilde{\mathbf{y}}_i^{*,r^{\text{OB}}} - \mathbb{Z}_i^{r^{\text{OB}}} \mathbf{b}_i^{r^{\text{OB}}}) \right\} \\ &\cdot \exp \left\{ -\frac{1}{2} (\mathbf{b}_i - \boldsymbol{\mu}^{(k)})^\top \Sigma^{-(k)} (\mathbf{b}_i - \boldsymbol{\mu}^{(k)}) \right\}, \end{aligned}$$

where $\tilde{\mathbf{y}}_i^{r^{\mathbb{N}}} = \mathbf{y}_i^{r^{\mathbb{N}}} - \mathbb{X}_i^{r^{\mathbb{N}}} \boldsymbol{\beta}^{r^{\mathbb{N}}}$ and $\tilde{\mathbf{y}}_i^{*,r^{\text{OB}}} = \mathbf{y}_i^{*,r^{\text{OB}}} - \mathbb{X}_i^{r^{\text{OB}}} \boldsymbol{\beta}^{r^{\text{OB}}}$. Constructing

$$\tilde{\mathbf{y}}_i = \begin{pmatrix} \vdots \\ \sqrt{\tau_{r^{\mathbb{N}}}^{(k)}} \tilde{\mathbf{y}}_i^{r^{\mathbb{N}}} \\ \vdots \\ \tilde{\mathbf{y}}_i^{*,r^{\text{OB}}} \\ \vdots \end{pmatrix}, \quad \begin{matrix} r^{\mathbb{N}} \in \mathcal{R}^{\text{Num}}, \\ r^{\text{OB}} \in \mathcal{R}^{\text{OB}}, \end{matrix} \quad \tilde{\mathbb{Z}}_i = \begin{pmatrix} \ddots & & & & \\ & \sqrt{\tau_{r^{\mathbb{N}}}^{(k)}} \mathbb{Z}_i^{r^{\mathbb{N}}} & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & \mathbb{Z}_i^{r^{\text{OB}}} \\ & & & & & \ddots \end{pmatrix}$$

we can simplify the above to

$$\exp \left\{ -\frac{1}{2} \left(\tilde{\mathbf{y}}_i - \tilde{\mathbb{Z}}_i \mathbf{b}_i \right)^\top \left(\tilde{\mathbf{y}}_i - \tilde{\mathbb{Z}}_i \mathbf{b}_i \right) - \frac{1}{2} \left(\mathbf{b}_i - \boldsymbol{\mu}^{(k)} \right)^\top \boldsymbol{\Sigma}^{- (k)} \left(\mathbf{b}_i - \boldsymbol{\mu}^{(k)} \right) \right\},$$

which after several algebraic operations and ignoring multiplicative constants becomes

$$\exp \left\{ -\frac{1}{2} \left(\tilde{\mathbf{b}}_i - \mathbf{b}_i \right)^\top \left[\tilde{\mathbb{Z}}_i^\top \tilde{\mathbb{Z}}_i + \boldsymbol{\Sigma}^{- (k)} \right] \left(\tilde{\mathbf{b}}_i - \mathbf{b}_i \right) \right\}, \quad (53)$$

where

$$\tilde{\mathbf{b}}_i = \left[\tilde{\mathbb{Z}}_i^\top \tilde{\mathbb{Z}}_i + \boldsymbol{\Sigma}^{- (k)} \right]^{-1} \left(\tilde{\mathbb{Z}}_i^\top \tilde{\mathbf{y}}_i + \boldsymbol{\Sigma}^{- (k)} \boldsymbol{\mu}^{(k)} \right).$$

Therefore, the full-conditioned distribution of \mathbf{b}_i for subject belonging to group $k = 1, \dots, K$ is

$$\mathbf{b}_i \left| \mathbb{Y}_i^{\text{N}}, \mathcal{C}_i; \mathbb{Y}_i^{\text{*,OB}}, U_i; \boldsymbol{\beta}, \boldsymbol{\tau}, \boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1} \sim N_{d_{\text{R}}} \left(\tilde{\mathbf{b}}_i, \left[\tilde{\mathbb{Z}}_i^\top \tilde{\mathbb{Z}}_i + \boldsymbol{\Sigma}^{- (k)} \right]^{-1} \right). \quad (54)$$

References

- 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, DOI 10.1214/08-AOAS205, URL <https://doi.org/10.1214/08-AOAS205>
- Albert JH, Chib S (1993) Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association* 88(422):669–679, DOI 10.2307/2290350
- Banfield D J, Raftery E A (1993) Model-based Gaussian and non-Gaussian clustering. *Biometrics* 49(3):803–821, URL <http://www.jstor.org/stable/2532201>
- Brooks S, Gelman A, Jones G, Meng X (2011) *Handbook for Markov chain Monte Carlo*, 2nd edn. Taylor & Francis
- Bruckers L, Molenberghs G, Drinkenburg P, Geys H (2016) A clustering algorithm for multivariate longitudinal data. *Journal of Biopharmaceutical Statistics* 26(4):725–741
- Celeux G, Martin O, Lavergne C (2005) Mixture of linear mixed models for clustering gene expression profiles from repeated microarray experiments. *Statistical Modelling* 5(3):243–267, DOI 10.1191/1471082X05st096oa
- De la Cruz-Mesía R, Quintana FA, Marshall G (2008) Model-based clustering for longitudinal data. *Computational Statistics and Data Analysis* 52(3):1441–1457, DOI 10.1016/j.csda.2007.04.005
- Dempster AP, Laird NM, Rubin DB (1977) Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society Series B (Methodological)* 39(1):1–38
- Fieuws S, Verbeke G (2004) Joint modelling of multivariate longitudinal profiles: Pitfalls of the random-effects approach. *Statistics in medicine* 23:3093–3104, DOI 10.1002/sim.1885
- Fieuws S, Verbeke G (2006) Pairwise fitting of mixed models for the joint modeling of multivariate longitudinal profiles. *Biometrics* 62(2):424–431
- Fraley C, Raftery AE (2002) Model-based clustering, discriminant analysis, and density estimation. *Journal of the American Statistical Association* 97(458):611–631, DOI 10.1198/016214502760047131
- Frühwirth-Schnatter S (2011) Panel data analysis: a survey on model-based clustering of time series. *Advances in Data Analysis and Classification* 5(4):251–280, DOI 10.1007/s11634-011-0100-0
- Geman S, Geman D (1984) Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence PAMI-6(6):721–741*, DOI 10.1109/TPAMI.1984.4767596
- Genz A (1992) Numerical computation of multivariate normal probabilities. *Journal of Computational and Graphical Statistics* 1(2):141–149
- Genz A, Bretz F, Miwa T, Mi X, Leisch F, Scheipl F, Hothorn T (2019) mvtnorm: Multivariate Normal and t Distributions. URL <https://CRAN.R-project.org/package=mvtnorm>, r package version 1.0-11
- 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, DOI 10.18637/jss.v028.i04
- Hastie T, Tibshirani R, Friedman J (2009) *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Second edn. Springer Science+Business Media, New York

- James GM, Sugar CA (2003) Clustering for sparsely sampled functional data. *Journal of the American Statistical Association* 98(462):397–408, DOI 10.1198/016214503000189
- Komárek A, Komárková L (2013) Clustering for multivariate continuous and discrete longitudinal data. *The Annals of Applied Statistics* 7(1):177–200, DOI 10.1214/12-AOAS580
- Komárek A, Komárková L (2014) Capabilities of R package mixAK for clustering based on multivariate continuous and discrete longitudinal data. *Journal of Statistical Software* 59(12):1–38, DOI 10.18637/jss.v059.i12
- Laird NM, Ware JH (1982) Random-effects models for longitudinal data. *Biometrics* 38(4):963–974
- Liu X, Yang MCK (2009) Simultaneous curve registration and clustering for functional data. *Computational Statistics and Data Analysis* 53(4):1361–1376, DOI 10.1016/j.csda.2008.11.019
- Ma P, Castillo-Davis CI, Zhong W, Liu JS (2006) A data-driven clustering method for time course gene expression data. *Nucleic Acids Research* 34(4):1261–1269, DOI 10.1093/nar/gkl013
- McNicholas PD, Murphy TB (2010) Model-based clustering of longitudinal data. *The Canadian Journal of Statistics* 38(1):153–168, DOI 10.1002/cjs.10047
- Molenberghs G, Verbeke G (2005) *Models for Discrete Longitudinal Data*. Springer, New York
- Neal RM (2000) Markov chain sampling methods for dirichlet process mixture models. *Journal of Computational and Graphical Statistics* 9(2):249–265
- Proust-Lima C, Philipps V, Diakite A, Lique B (2017) Estimation of extended mixed models using latent classes and latent processes: The R package lcmm. *Journal of Statistical Software* 78(2):1–56, DOI doi:10.18637/jss.v078.i02
- Raftery AE, Dean N (2006) Variable selection for model-based clustering. *Journal of the American Statistical Association* 101(473):168–178
- Robert CP (2001) *The Bayesian Choice*, 2nd edn. Springer-Verlag New York
- Stephens M (2000) Dealing with label switching in mixture models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 62(4):795–809
- Tanner MA, Wong WH (1987) The calculation of posterior distributions by data augmentation. *Journal of the American Statistical Association* 82(398):528–550, DOI 10.2307/2289457
- Verbeke G, Lesaffre E (1996) A linear mixed-effects model with heterogeneity in the random-effects population. *Journal of the American Statistical Association* 91(433):217–221, DOI 10.1080/01621459.1996.10476679
- Villarroel L, Marshall G, Barón AE (2009) Cluster analysis using multivariate mixed effects models. *Statistics in Medicine* 28(20):2552–2565, DOI 10.1002/sim.3632