



FACULTY
OF MATHEMATICS
AND PHYSICS
Charles University

Arnošt Komárek & Jan Vávra

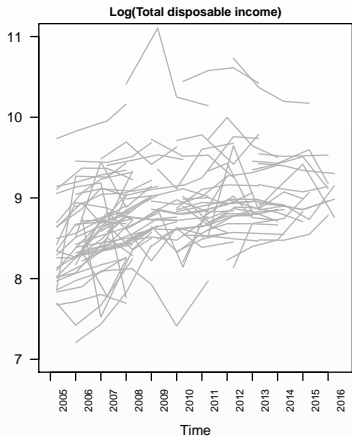
**Classification Based on Mixed Type
Numeric, Ordinal and Binary
Longitudinal Data**

I.

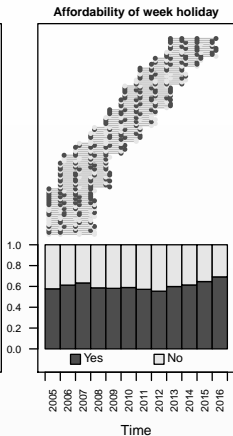
Introduction

(Multivariate and mixed type) longitudinal data

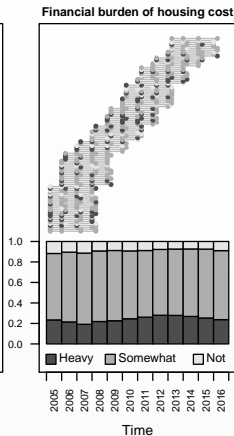
Numeric



Binary



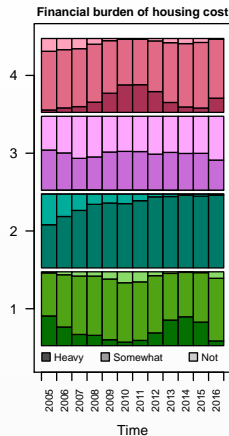
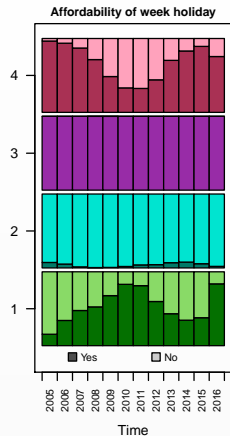
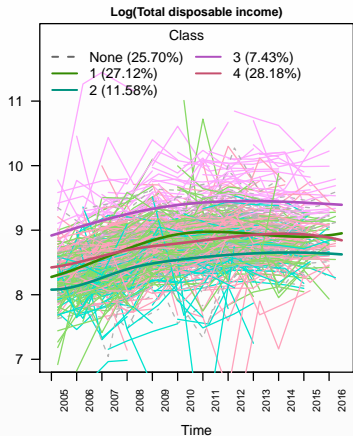
Ordinal



$$\mathbf{Y}_i \equiv \underbrace{(Y_{i,1}^1, \dots, Y_{i,n_i}^1)}_{\mathbf{Y}_i^1}, \underbrace{(Y_{i,1}^2, \dots, Y_{i,n_i}^2)}_{\mathbf{Y}_i^2}, \underbrace{(Y_{i,1}^3, \dots, Y_{i,n_i}^3)}_{\mathbf{Y}_i^3}, \quad i = 1, \dots, n$$

$$\mathbf{v}_{i,j}^r \equiv (t_{i,j}, \dots): \text{additional covariates}, \quad j = 1, \dots, n_i, \quad r = 1, \dots, R(=3)$$

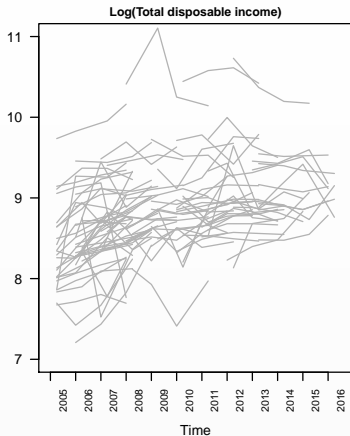
Clustering (unsupervised classification)



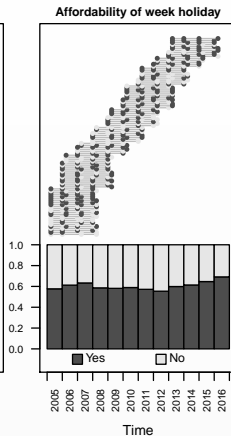
- “Classical” data: $\mathbf{Y}_1, \dots, \mathbf{Y}_n \in \mathbb{R}^p$
- K-means, hierarchical methods, ...
 - based on a **distance** in p -dimensional Euclidean space
- (Mixed type) longitudinal data
 - $\mathbf{Y}_i \in \mathbb{R}^{R n_i}$ – different numbers of measurements per subject
 - irregularly spaced in time
 - some $Y_{i,j} \in \{0, 1\}, \in \{0, 1, 2\}, \dots$
 - distance?

(Multivariate and mixed type) longitudinal data

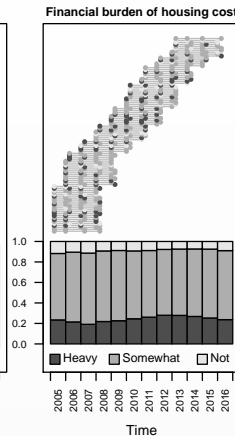
Numeric



Binary



Ordinal



$$Y_i \equiv \underbrace{(Y_{i,1}^1, \dots, Y_{i,n_i}^1)}_{Y_i^1}, \underbrace{(Y_{i,1}^2, \dots, Y_{i,n_i}^2)}_{Y_i^2}, \underbrace{(Y_{i,1}^3, \dots, Y_{i,n_i}^3)}_{Y_i^3}, \quad i = 1, \dots, n$$

$$\mathbf{v}_{i,j}^r \equiv (t_{i,j}, \dots): \text{additional covariates}, \quad j = 1, \dots, n_i, \quad r = 1, \dots, R(=3)$$

II.

Model based clustering

Data for subject i , \mathcal{D}_i :

$$\mathbf{Y}_i \equiv \left(\underbrace{Y_{i,1}^1, \dots, Y_{i,n_i}^1}_{\mathbf{Y}_i^1}, \underbrace{Y_{i,1}^2, \dots, Y_{i,n_i}^2}_{\mathbf{Y}_i^2}, \underbrace{Y_{i,1}^3, \dots, Y_{i,n_i}^3}_{\mathbf{Y}_i^3} \right), \quad i = 1, \dots, n$$

$\mathbf{v}_{i,j}^r \equiv (t_{i,j}, \dots)$: additional covariates, $j = 1, \dots, n_i$, $r = 1, \dots, R (= 3)$
 $\rightarrow \mathcal{C}_i$

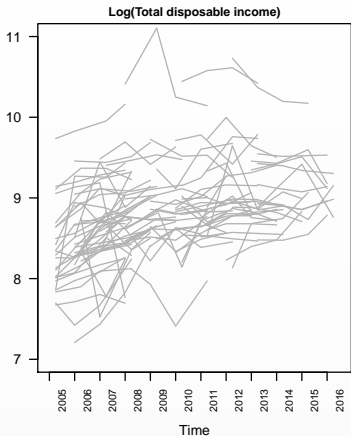
- (Whole) population consists of K subpopulations.
In advance, we do not know **who is who**.
- Data in subpopulation $k \in \{1, \dots, K\}$ follow certain **statistical** model

$\rightarrow f_k(\mathbf{y}_i; \boldsymbol{\xi}^k, \boldsymbol{\xi}, \mathcal{C}_i)$

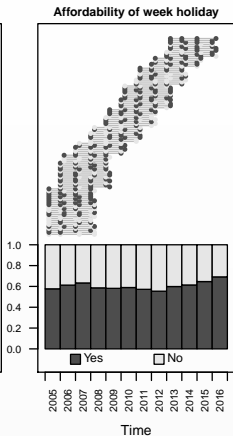
- $\boldsymbol{\xi}^k$: group-specific (unknown) model parameters
- $\boldsymbol{\xi}$: (unknown) model parameters shared by all groups

(Multivariate and mixed type) longitudinal data

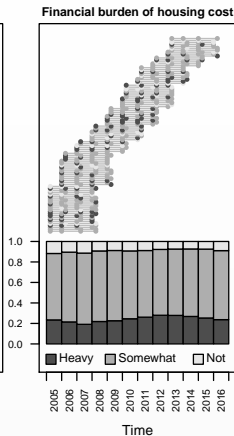
Numeric



Binary



Ordinal



$$\mathbf{Y}_i \equiv \underbrace{(Y_{i,1}^1, \dots, Y_{i,n_i}^1)}_{\mathbf{Y}_i^1}, \underbrace{(Y_{i,1}^2, \dots, Y_{i,n_i}^2)}_{\mathbf{Y}_i^2}, \underbrace{(Y_{i,1}^3, \dots, Y_{i,n_i}^3)}_{\mathbf{Y}_i^3}, \quad i = 1, \dots, n$$

$$\mathbf{v}_{i,j}^r \equiv (t_{i,j}, \dots): \text{additional covariates}, \quad j = 1, \dots, n_i, \quad r = 1, \dots, R(=3)$$

Simple example

- Group-specific model \equiv “simple” gaussian linear model

$$Y_{i,j} = \beta_0^k + \beta_1^k t_{i,j} + \varepsilon_{i,j}, \quad \varepsilon_{i,j} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$$

- $\xi^k \equiv \beta^k \equiv (\beta_0^k, \beta_1^k)$

- $\xi \equiv \sigma^2$

- $f_k(\mathbf{y}_i; \xi^k, \xi, C_i) = \prod_{j=1}^{n_i} \varphi(y_{i,j}; \beta_0^k + \beta_1^k t_{i,j}, \sigma^2)$

- Group-specific model \equiv **linear mixed model**

$$Y_{i,j} = \beta_0^k + \beta_1^k t_{i,j} + \mathbf{b}_{i,0} + \mathbf{b}_{i,1} t_{i,j} + \varepsilon_{i,j}, \quad \varepsilon_{i,j} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\mathbf{b}_i = (\mathbf{b}_{i,0}, \mathbf{b}_{i,1})^\top \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbb{D})$$

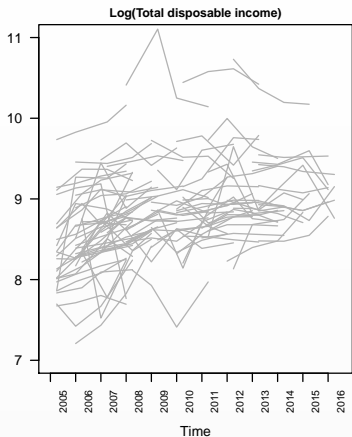
- $\xi^k \equiv \beta^k \equiv (\beta_0^k, \beta_1^k)$
- $\xi \equiv (\sigma^2, \Sigma)$
- $f_k(\mathbf{y}_i; \xi^k, \xi, C_i) = \varphi(\mathbf{y}_i; \mathbb{X}_i \beta^k, \mathbb{V}_i),$

$$\mathbb{X}_i = \begin{pmatrix} 1 & t_{i,1} \\ \vdots & \vdots \\ 1 & t_{i,n_i} \end{pmatrix}$$

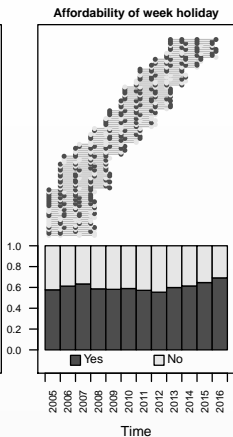
$$\mathbb{V}_i = \mathbb{X}_i \Sigma \mathbb{X}_i^\top + \sigma^2 \mathbf{I}_{n_i}$$

(Multivariate and mixed type) longitudinal data

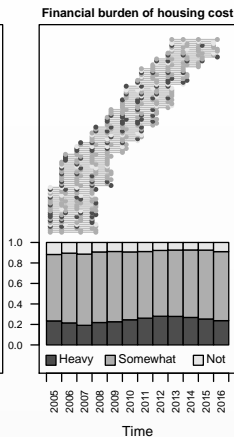
Numeric



Binary



Ordinal



$$\mathbf{Y}_i \equiv \underbrace{(Y_{i,1}^1, \dots, Y_{i,n_i}^1)}_{\mathbf{Y}_i^1}, \underbrace{(Y_{i,1}^2, \dots, Y_{i,n_i}^2)}_{\mathbf{Y}_i^2}, \underbrace{(Y_{i,1}^3, \dots, Y_{i,n_i}^3)}_{\mathbf{Y}_i^3}, \quad i = 1, \dots, n$$

$$\mathbf{v}_{i,j}^r \equiv (t_{i,j}, \dots): \text{additional covariates}, \quad j = 1, \dots, n_i, \quad r = 1, \dots, R(=3)$$

Model based clustering

- Data in subpopulation $k \in \{1, \dots, K\}$ follow certain **statistical** model

$$\rightarrow f_k(\mathbf{y}_i; \xi^k, \xi, C_i)$$

- (Yet) unknown allocation of subject i in either of K groups:
 $U_i \in \{1, \dots, K\}$

$$f_k(\mathbf{y}_i; \xi^k, \xi, C_i) = p(\mathbf{y}_i | U_i = k; \xi^k, \xi, C_i)$$

- (Unknown) proportions of K groups in the whole population:

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

$$\mathbf{w} = (w_1, \dots, w_K)^\top \in (0, 1)^K$$

- Above implies the (marginal) distribution of response

$$\begin{aligned} f(\mathbf{y}_i; \boldsymbol{\theta}, C_i) &= p(\mathbf{y}_i; \boldsymbol{\theta}, C_i) \\ &= \sum_{k=1}^K p(\mathbf{y}_i | U_i = k; \boldsymbol{\xi}^k, \boldsymbol{\xi}, C_i) P(U_i = k; \mathbf{w}) \\ &= \underbrace{\sum_{k=1}^K w_k f_k(\mathbf{y}_i; \boldsymbol{\xi}^k, \boldsymbol{\xi}, C_i)}_{\text{mixture of models}} \end{aligned}$$

- Unknown parameters of (the whole) model:

$$\boldsymbol{\theta} \equiv (\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^K)$$

and a mixture likelihood (if independence across subjects assumed)

$$L(\boldsymbol{\theta}) = \prod_{i=1}^n f(\mathbf{y}_i; \boldsymbol{\theta}, C_i)$$

Classification?

- (Not yet) really Bayesian analysis but a simple use of the Bayes theorem ($k = 1, \dots, K$):

$$\begin{aligned} u_{i,k}(\boldsymbol{\theta}) &:= P(U_i = k \mid \mathbf{Y}_i = \mathbf{y}_i; \boldsymbol{\theta}, C_i) \\ &= \frac{p(\mathbf{y}_i \mid U_i = k; \boldsymbol{\xi}^k, \boldsymbol{\xi}, C_i) P(U_i = k; \mathbf{w})}{\sum_{l=1}^K p(\mathbf{y}_i \mid U_i = l; \boldsymbol{\xi}^l, \boldsymbol{\xi}, C_i) P(U_i = l; \mathbf{w})} \\ &= \frac{w_k f_k(\mathbf{y}_i; \boldsymbol{\xi}^k, \boldsymbol{\xi}, C_i)}{\sum_{l=1}^K w_l f_l(\mathbf{y}_i; \boldsymbol{\xi}^l, \boldsymbol{\xi}, C_i)} \end{aligned}$$

Classification?

- If we knew all model parameters $\theta \equiv (\mathbf{w}, \xi, \xi^1, \dots, \xi^K)$, possible classification rule (“prediction” of the allocation variable U_i):

$$\hat{U}_i = \operatorname{argmax}_{k=1, \dots, K} u_{i,k}(\theta)$$

- In case we do not know θ , just **estimate** it ($\longrightarrow \hat{\theta}$) and then

$$\hat{U}_i = \operatorname{argmax}_{k=1, \dots, K} u_{i,k}(\hat{\theta})$$

- Likelihood

$$L(\theta) = \prod_{i=1}^n \left\{ \sum_{k=1}^K w_k f_k(\mathbf{y}_i; \xi^k, \xi, C_i) \right\}$$

$$\stackrel{\text{LMM}}{=} \prod_{i=1}^n \left\{ \sum_{k=1}^K w_k \varphi(\mathbf{y}_i; \mathbb{X}_i \beta^k, \mathbb{X}_i \boldsymbol{\Sigma} \mathbb{X}_i^T + \sigma^2 \mathbf{I}_{n_i}) \right\}$$

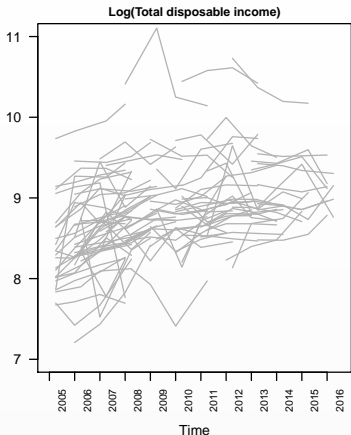
- $\hat{\theta} = \operatorname{argmax} L(\theta)$
- EM algorithm?

III.

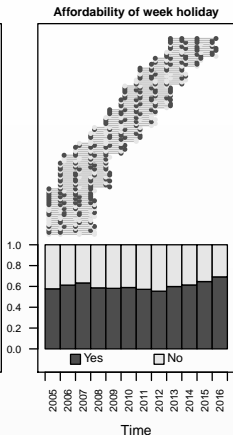
Model for mixed type longitudinal data

(Multivariate and mixed type) longitudinal data

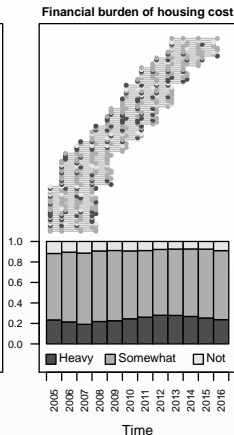
Numeric



Binary



Ordinal



$$Y_i \equiv \underbrace{(Y_{i,1}^1, \dots, Y_{i,n_i}^1)}_{Y_i^1}, \underbrace{(Y_{i,1}^2, \dots, Y_{i,n_i}^2)}_{Y_i^2}, \underbrace{(Y_{i,1}^3, \dots, Y_{i,n_i}^3)}_{Y_i^3}, \quad i = 1, \dots, n$$

$$\mathbf{v}_{i,j}^r \equiv (t_{i,j}, \dots): \text{additional covariates, } j = 1, \dots, n_i, r = 1, \dots, R(=3) \rightarrow C_i^r$$

Model for mixed type longitudinal data

- Hierarchical model
- **Numeric** outcome: linear mixed model (LME), i.e.,

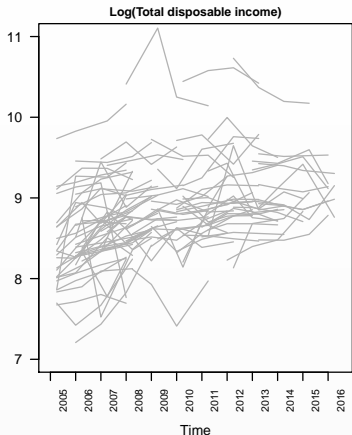
$$\mathbf{Y}_i^r \mid \mathbf{b}_i^r; \mathcal{C}_i^r \sim \mathcal{N}_{n_i}(\boldsymbol{\eta}_i^r, \tau_r^{-1} \mathbb{I}_{n_i}),$$

$$\boldsymbol{\eta}_i^r = \mathbb{X}_i^r \boldsymbol{\beta}^r + \mathbb{Z}_i^r \mathbf{b}_i^r$$

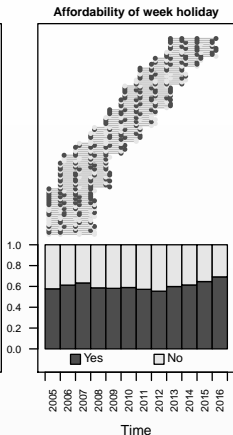
$$\mathbf{b}_i^r \stackrel{\text{i.i.d.}}{\sim} \text{wait a little bit}$$

(Multivariate and mixed type) longitudinal data

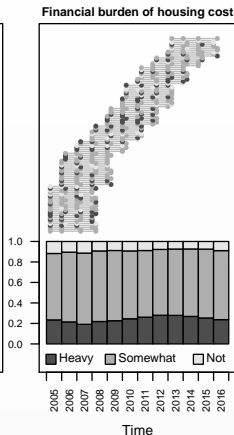
Numeric



Binary



Ordinal



$$Y_i \equiv \underbrace{(Y_{i,1}^1, \dots, Y_{i,n_i}^1)}_{Y_i^1}, \underbrace{(Y_{i,1}^2, \dots, Y_{i,n_i}^2)}_{Y_i^2}, \underbrace{(Y_{i,1}^3, \dots, Y_{i,n_i}^3)}_{Y_i^3}, \quad i = 1, \dots, n$$

$$\mathbf{v}_{i,j}^r \equiv (t_{i,j}, \dots): \text{additional covariates, } j = 1, \dots, n_i, r = 1, \dots, R(=3) \rightarrow C_i^r$$

Model for mixed type longitudinal data

- **Binary and ordinal** outcome: linear mixed model (LME) for a **latent** numeric variable
- **Thresholding** concept

$$\text{we observe } Y_{i,j}^r = l \iff \gamma_l^r < Y_{i,j}^{*,r} \leq \gamma_{l+1}^r$$

- $Y_{i,j}^{*,r}$: **latent** (unobservable) numeric variable
- $-\infty = \gamma_0^r < \gamma_1^r < \dots < \gamma_{L_r}^r = \infty$: unknown thresholds
- For identifiability purposes: $\gamma_1^r = \text{fixed const}$
- Then as before:

$$\mathbf{Y}_i^{*,r} \mid \mathbf{b}_i^r; C_i^r \sim \mathcal{N}_{n_i}(\boldsymbol{\eta}_i^r, \tau_r^{-1} \mathbb{I}_{n_i}),$$

$\tau_r = 1$ for identifiability purposes

$$\boldsymbol{\eta}_i^r = \mathbb{X}_i^r \boldsymbol{\beta}^r + \mathbb{Z}_i^r \mathbf{b}_i^r$$

- How to model dependencies across different outcomes on one subject?
- Joint distribution for all random effects

$$\mathbf{b}_i = \begin{pmatrix} \mathbf{b}_i^{\text{N}} \\ \mathbf{b}_i^{\text{OB}} \end{pmatrix} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_{d^{\text{R}}} \left(\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}^{\text{N}} \\ \boldsymbol{\mu}^{\text{OB}} \end{pmatrix}, \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}^{\text{N}} & \boldsymbol{\Sigma}^{\text{NOB}} \\ \boldsymbol{\Sigma}^{\text{OBN}} & \boldsymbol{\Sigma}^{\text{OB}} \end{pmatrix} \right)$$

- $\boldsymbol{\mu}$: (unknown) mean of random effects
- $\boldsymbol{\Sigma}$: (unknown) covariance matrix of random effects

$$\begin{aligned}
 & \rho(\mathbb{Y}_i \mid \boldsymbol{\beta}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\tau}, \boldsymbol{\gamma}; \mathcal{C}_i) \\
 &= \int \int \rho(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{\text{OB}}, \mathbb{Y}_i^{*,\text{OB}}, \mathbf{b}_i \mid \boldsymbol{\beta}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\tau}, \boldsymbol{\gamma}; \mathcal{C}_i) d\mathbf{b}_i d\mathbb{Y}_i^{*,\text{OB}} \\
 &= \int \int \underbrace{\rho(\mathbb{Y}_i^{\text{OB}} \mid \mathbb{Y}_i^{*,\text{OB}}, \boldsymbol{\gamma})}_{\text{thresholding}} \cdot \underbrace{\rho(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{*,\text{OB}} \mid \mathbf{b}_i, \boldsymbol{\beta}, \boldsymbol{\tau}; \mathcal{C}_i)}_{\text{MV LME}} \\
 &\quad \cdot \underbrace{\rho(\mathbf{b}_i \mid \boldsymbol{\mu}, \boldsymbol{\Sigma})}_{\text{normality}} d\mathbf{b}_i d\mathbb{Y}_i^{*,\text{OB}}
 \end{aligned}$$

- In principle, the likelihood corresponds to the CDF of **truncated multivariate** normal distribution

$$p(\mathbb{Y}_i^{\text{OB}} \mid \mathbb{Y}_i^{*,\text{OB}}, \gamma) = \prod_{r \in \mathcal{R}^{\text{OB}}} \prod_{j=1}^{n_i} \left[\sum_{l=0}^{L^r-1} \mathbb{I}_{\{l\}}(\mathbf{y}_{i,j}^r) \mathbb{I}_{(\gamma_l^r, \gamma_{l+1}^r]}(\mathbf{y}_i^{*,\text{OB}}) \right]$$

$$\begin{aligned} p(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{*,\text{OB}} \mid \mathbf{b}_i, \boldsymbol{\beta}, \boldsymbol{\tau}; \mathcal{C}_i) \\ = \prod_{r \in \mathcal{R}^{\text{Num}}} \prod_{j=1}^{n_i} \varphi(\mathbf{y}_{i,j}^r; \boldsymbol{\eta}_{i,j}^r, \boldsymbol{\tau}_r^{-1}) \cdot \prod_{r \in \mathcal{R}^{\text{OB}}} \prod_{j=1}^{n_i} \varphi(\mathbf{y}_{i,j}^{*,r}; \boldsymbol{\eta}_{i,j}^r, \mathbf{1}) \end{aligned}$$

$$p(\mathbf{b}_i \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \varphi(\mathbf{b}_i; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Mixture likelihood

- Do not forget, we will **mix** the above model
- On top of random effects and latent variables for binary/ordinal outcomes, there will be also latent component allocations $U_i \dots$
- The **mixture likelihood** is then

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

$$\begin{aligned} p\left(\mathbb{Y}_i \mid \beta, \mu, \Sigma, \tau, \gamma; C_i\right) &= \int \int \underbrace{p\left(\mathbb{Y}_i^{\text{OB}} \mid \mathbb{Y}_i^{*,\text{OB}}, \gamma\right)}_{\text{thresholding}} \cdot \underbrace{p\left(\mathbb{Y}_i^{\text{N}}, \mathbb{Y}_i^{*,\text{OB}} \mid \mathbf{b}_i, \beta, \tau; C_i\right)}_{\text{MV LME}} \\ &\quad \cdot \underbrace{p\left(\mathbf{b}_i \mid \mu, \Sigma\right)}_{\text{normality}} d\mathbf{b}_i d\mathbb{Y}_i^{*,\text{OB}} \end{aligned}$$

Maximum likelihood?

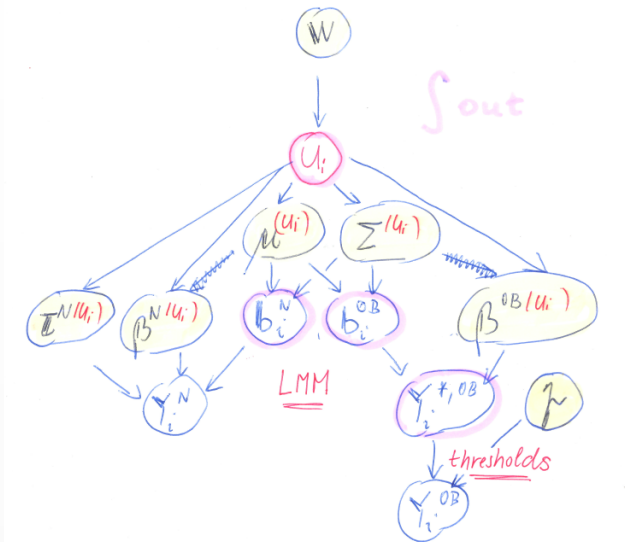
- Maximum-likelihood estimation?
- EM algorithm?
- Good luck. . .



IV.

Bayesian inference

Hierarchical model

The model is fully **hierarchical**:



- Specify a prior distribution for θ (in a *reasonable* way)
- Principles of Bayesian Data Augmentation (BDA) can be used towards Gibbs sampling and MCMC to obtain the posterior sample $\theta^{(m)}$, $m = 1, \dots, M$ for inference
- Has been implemented (by JV) in  and 

Bayesian classification

Remember ($k = 1, \dots, K$):

$$\begin{aligned}u_{i,k}(\boldsymbol{\theta}) &= \mathbf{P}(U_i = k \mid \mathbf{Y}_i = \mathbf{y}_i; \boldsymbol{\theta}, C_i) \\ &= \frac{w_k f_k(\mathbf{y}_i; \boldsymbol{\xi}^k, \boldsymbol{\xi}, C_i)}{\sum_{l=1}^K w_l f_l(\mathbf{y}_i; \boldsymbol{\xi}^l, \boldsymbol{\xi}, C_i)}\end{aligned}$$

Here:

$$\begin{aligned}f_k(\mathbf{y}_i; \boldsymbol{\xi}^k, \boldsymbol{\xi}, C_i) &\equiv p(\mathbf{Y}_i \mid \boldsymbol{\beta}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)}, \boldsymbol{\tau}^{(k)}, \boldsymbol{\gamma}; C_i) \\ &= \int \int \underbrace{p(\mathbf{Y}_i^{\text{OB}} \mid \mathbf{Y}_i^{*,\text{OB}}, \boldsymbol{\gamma})}_{\text{thresholding}} \cdot \underbrace{p(\mathbf{Y}_i^{\text{N}}, \mathbf{Y}_i^{*,\text{OB}} \mid \mathbf{b}_i, \boldsymbol{\beta}^{(k)}, \boldsymbol{\tau}^{(k)}; C_i)}_{\text{MV LME}} \\ &\quad \cdot \underbrace{p(\mathbf{b}_i \mid \boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)})}_{\text{normality}} d\mathbf{b}_i d\mathbf{Y}_i^{*,\text{OB}}\end{aligned}$$

Bayesian classification

- For each i and k , quantity $u_{i,k}(\boldsymbol{\theta})$ is just some function of unknown parameters
- In a Bayesian classification procedure, rather than $u_{i,k}(\hat{\boldsymbol{\theta}})$, one would use some estimate of $u_{i,k}(\boldsymbol{\theta})$ directly, e.g.,
- The **posterior mean**

$$\hat{U}_{i,k} = \int u_{i,k}(\boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathcal{D}) d\boldsymbol{\theta} \approx \frac{1}{M} \sum_{m=1}^M u_{i,k}(\boldsymbol{\theta}^{(m)})$$

- “Small” complication: integrals over the multivariate normal density must be calculated
- Next to the point estimate, also the (95% HPD) **credible interval** for $u_{i,k}(\boldsymbol{\theta})$ can be calculated to quantify uncertainty

$$\longrightarrow (\hat{U}_{i,k}^{LOW}, \hat{U}_{i,k}^{UPP})$$

Bayesian classification

- Traditional rule: $\hat{U}_i = \operatorname{argmax}_{k=1,\dots,K} \hat{U}_{i,k}$

→ may lead to high misclassification rates if two (or more) clusters are not really apart

- Uncertainty in classification may, e.g., be taken into account as follows:

Subject i is classified in group $\hat{U}_i = k$ if and only if the lower bound $\hat{U}_{i,k}^{\text{LOW}}$ of the credible interval is higher than any other upper bound $\hat{U}_{i,l}^{\text{UPP}}$, $l \neq k$. If this does not happen, subject i remains unclassified.

- fills clusters with their most typical representatives
- keep indecisive subjects aside

Classification of a “new” subject

- The whole procedure can also be used to classify a “new” subject without re-fitting the model (running the MCMC procedure)
- Allocation probabilities are just a function of model parameters

$$\begin{aligned} u_{new,k}(\theta) &= P(U_{new} = k \mid \mathbf{Y}_{new} = \mathbf{y}_{new}; \theta, C_{new}) \\ &= \frac{w_k f_k(\mathbf{y}_{new}; \xi^k, \xi, C_{new})}{\sum_{l=1}^K w_l f_l(\mathbf{y}_{new}; \xi^l, \xi, C_{new})} \end{aligned}$$

- The posterior mean (and the credible interval) does not need to include the “new” observation in the posterior distribution

$$\hat{U}_{new,k} = \int u_{new,k}(\theta) p(\theta \mid \mathcal{D}(old)) d\theta \approx \frac{1}{M} \sum_{m=1}^M u_{new,k}(\theta^{(m)})$$

Classification of a “new” subject

- Allocation probabilities are just a function of model parameters

$$\begin{aligned} u_{new,k}(\theta) &= P(U_{new} = k \mid \mathbf{Y}_{new} = \mathbf{y}_{new}; \theta, C_{new}) \\ &= \frac{w_k f_k(\mathbf{y}_{new}; \xi^k, \xi, C_{new})}{\sum_{l=1}^K w_l f_l(\mathbf{y}_{new}; \xi^l, \xi, C_{new})} \end{aligned}$$

- With longitudinal data $\mathbf{Y}_{new} = (Y_{new,1}, Y_{new,2}, \dots, Y_{new,n_{new}})^\top$, classification can proceed **dynamically**. We will follow the “new” subject only until the moment when it is clear to which cluster it belongs to

V.

Simulation

Simulation setting

- Three longitudinal outcomes: numeric, binary, ordinal (3 levels)
- Four measurement occasions: $0 < t_{i,1} < t_{i,2} < t_{i,3} < t_{i,4} < 1$
- Linear predictor (for all outcomes)

$$1 X_{i,j}^1 - 2 X_{i,j}^2 + \dots,$$

$$X_{i,j}^1 \stackrel{\text{i.i.d.}}{\sim} \text{Alt}(0.5), \quad X_{i,j}^2 \stackrel{\text{i.i.d.}}{\sim} \text{Unif}(0, 1)$$

- ...

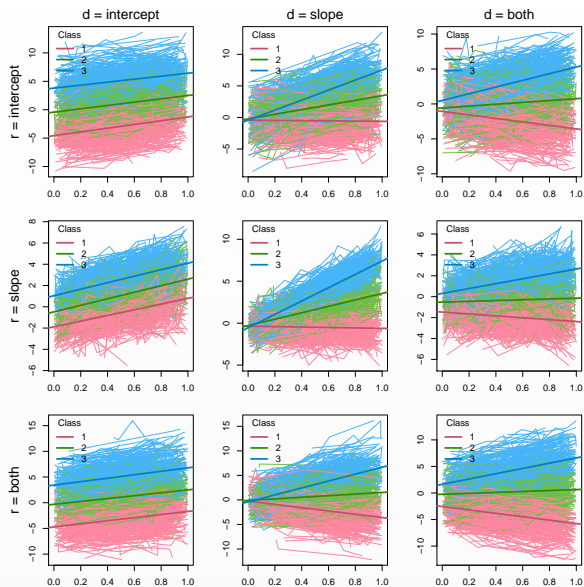
- ① (r = intercept): $b_{0,i} + \beta_1 t_{i,j}$, $E(b_{0,i}) = \beta_0$
- ② (r = slope): $\beta_0 + b_{1,i} t_{i,j}$, $E(b_{1,i}) = \beta_1$
- ③ (r = both): $b_{0,i} + b_{1,i} t_{i,j}$, $E(b_{0,i}, b_{1,i}) = (\beta_0, \beta_1)$

- Class specific?

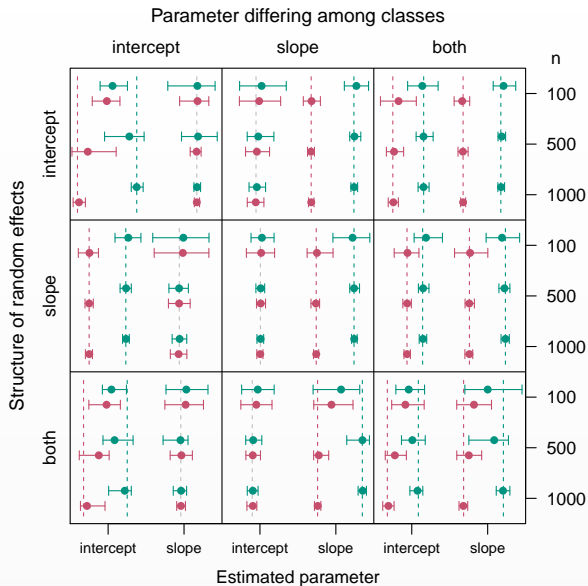
- ① (d = intercept): $\beta_0 \equiv \beta_0^{(k)}$
- ② (d = slope): $\beta_1 \equiv \beta_1^{(k)}$
- ③ (d = both): $\beta_0 \equiv \beta_0^{(k)}, \beta_1 \equiv \beta_1^{(k)}$

→ $3 \times 3 = 9$ scenarios with $K = 2, 3$ and $n = 100, 500, 1000$

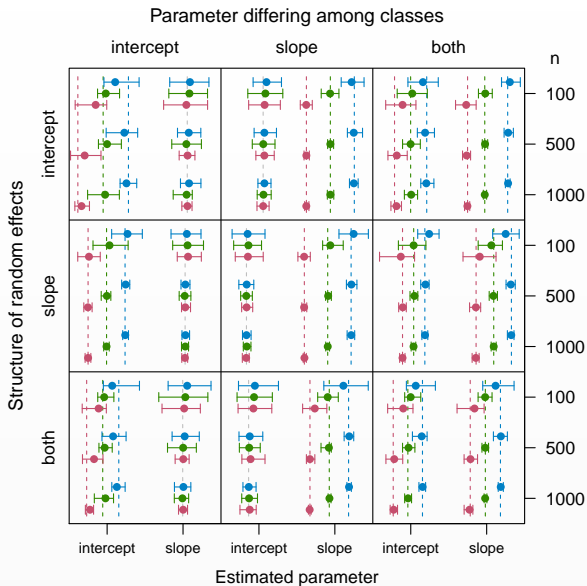
Simulation setting, $K = 3$



Statistical properties of posterior means, $K = 2$



Statistical properties of posterior means, $K = 3$



Classification ability

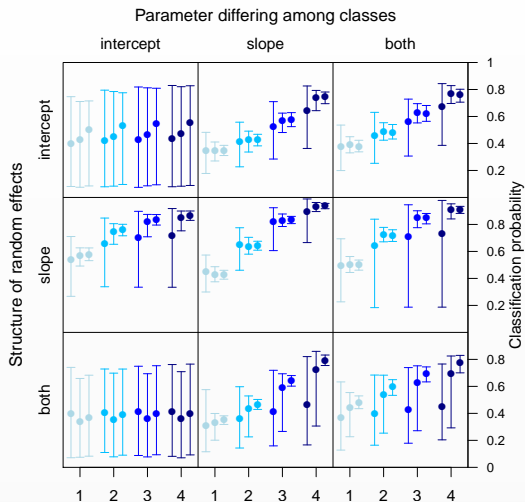
r	d	n	K = 2			K = 3		
			Correct [%]	Uncl. [%]	Miscl. [%]	Correct [%]	Uncl. [%]	Miscl. [%]
intercept	intercept	100	27.0 (17.2)	63.2 (25.4)	9.8 (13.7)	23.0 (17.5)	70.2 (21.0)	6.8 (9.4)
		500	62.5 (27.2)	33.0 (27.3)	4.4 (3.8)	44.3 (20.6)	50.8 (22.1)	4.9 (4.4)
		1000	85.1 (6.7)	10.1 (7.1)	4.8 (0.9)	58.6 (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)

Classification ability

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)

Dynamic classification

- Remember: $u_{i,k}(\theta) = P(U_i = k \mid \mathbf{Y}_i = \mathbf{y}_i; \theta, C_i)$
- $\mathbf{Y}_i = Y_{i,1}$, then $\mathbf{Y}_i = (Y_{i,1}, Y_{i,2})$, then $\mathbf{Y}_i = (Y_{i,1}, Y_{i,2}, Y_{i,3}), \dots$

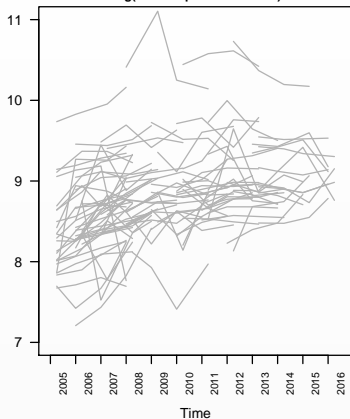


VI.

EU-SILC (Czech)

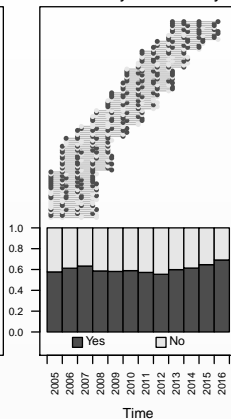
Numeric

Log(Total disposable income)



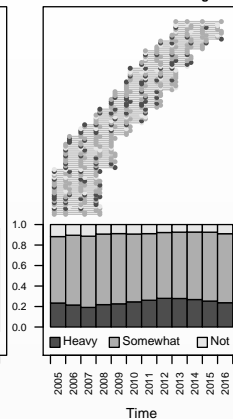
Binary

Affordability of week holiday



Ordinal

Financial burden of housing cost



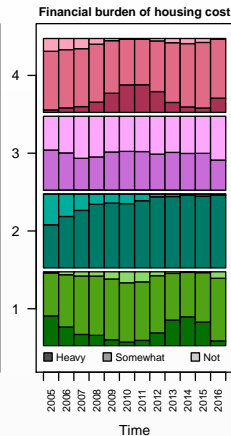
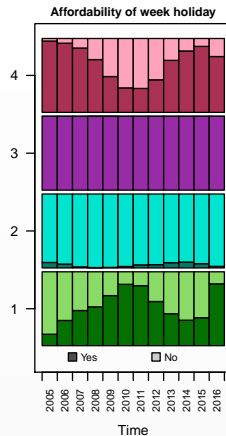
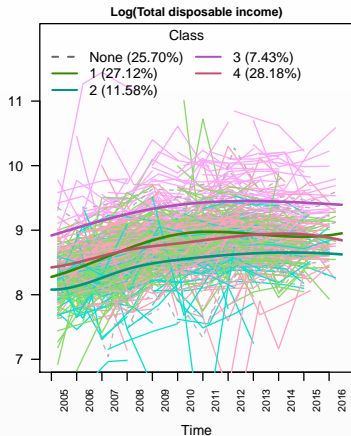
Model

- Linear predictor (for outcome $r \in \{1, 2, 3\}$):

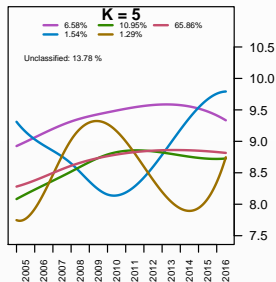
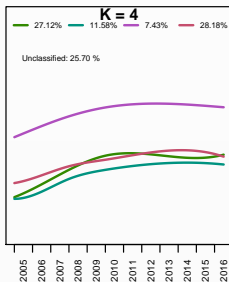
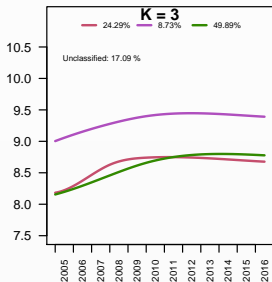
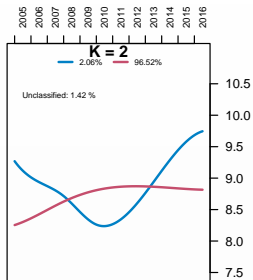
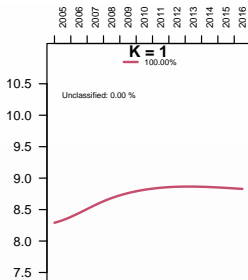
$$b_{0,i}^r + \underbrace{\beta_0^r + \beta_1^r B_1(t_{i,j}) + \cdots + \beta_5^r B_5(t_{i,j})}_{\text{spline in time}} + \beta_6^r \underbrace{S_{i,j}}_{\text{weighted family size}}$$

\implies all β 's possibly class specific (depend on k)

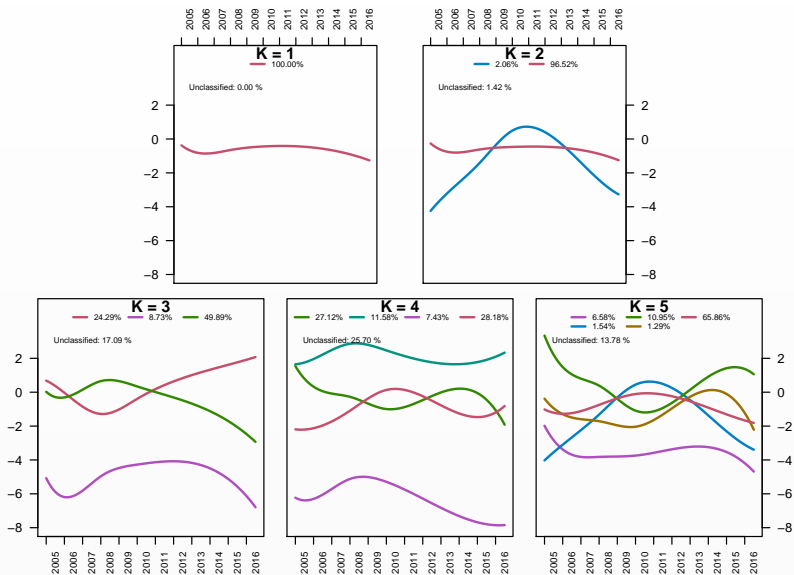
- Random effect vector: $\mathbf{b}_{0,i} = (b_{0,i}^1, b_{0,i}^2, b_{0,i}^3)^\top$



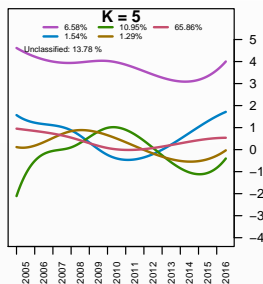
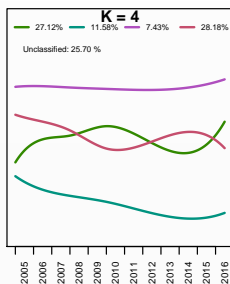
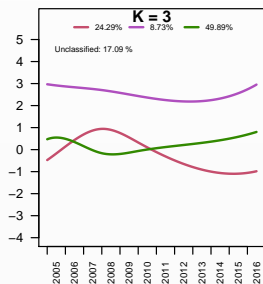
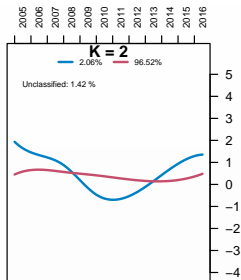
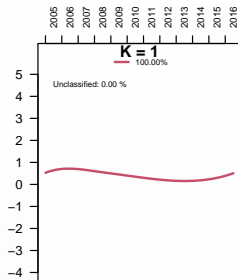
EU-SILC (Czech), log(total disposable income)

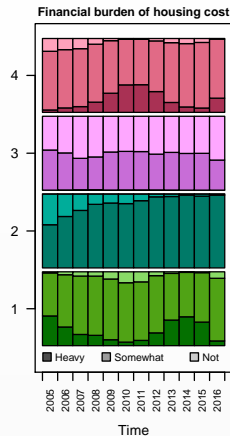
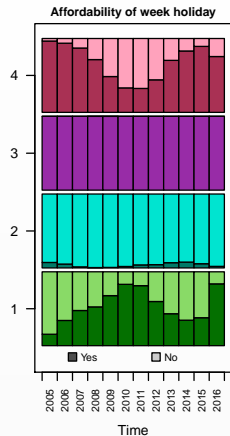
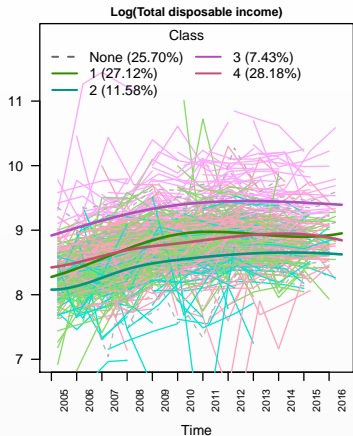


EU-SILC (Czech), affordability of week holiday

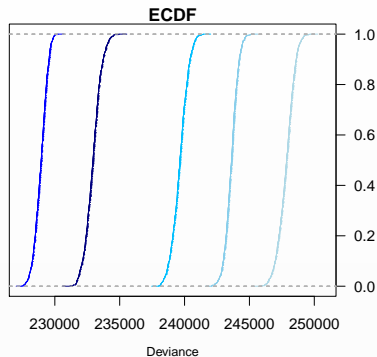
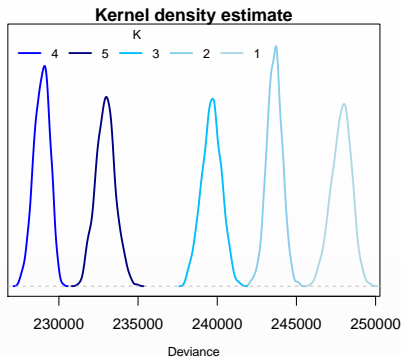


EU-SILC (Czech), financial burden of housing cost)





How many groups?



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

VII.

Conclusions

Conclusions

- Model based clustering is a general clustering method applicable with almost any data structures
- The only thing we need to do is to specify a **model** for data at hand which also expresses presumable differences between clusters
- Bayesian MCMC calculation solves (relatively easily) problems with a mixture likelihood and other possibly ugly integrals as soon as the model can be written hierarchically
- With MCMC based inference, also uncertainty in classification can (realively easily) be taken into account
- MBC also naturally allows for development of, in fact, a **discriminant** procedure allowing to classify “new” subjects (diagnosis in medicine, . . .)

- Also **nominal** categorical data are quite frequent (especially in social sciences. . .)
 - Thresholding approach could in principle be used as well to link categorical outcome to the numeric one
- Interpretation of the model parameters (trends) would, however, be somehow strange
- Alternative: sort of **generalized linear mixed model (GLMM)** to be used for categorical outcomes

**Vielen Dank für Ihre
Aufmerksamkeit!**