



FACULTY
OF MATHEMATICS
AND PHYSICS
Charles University

Jan Vávra, Arnošt Komárek


Department of Probability and Mathematical Statistics

**Classification Based on Longitudinal Data of
a Mixed Type**

16th July 2019

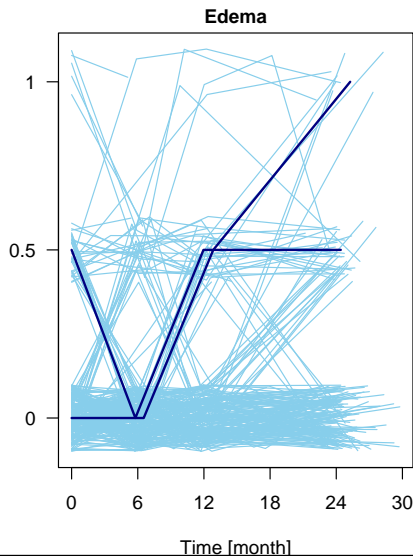
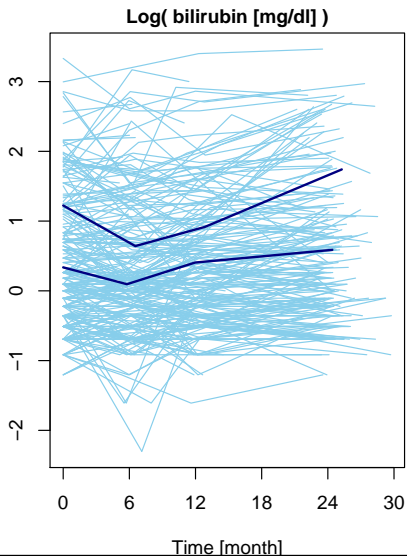
Mayo Clinic PBC Data (1974 – 1984)

Flemming and Harrington (1991)

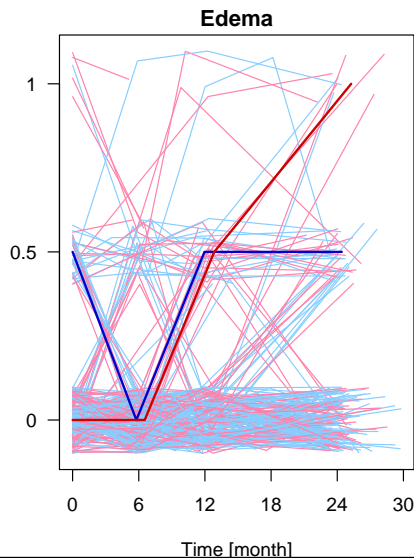
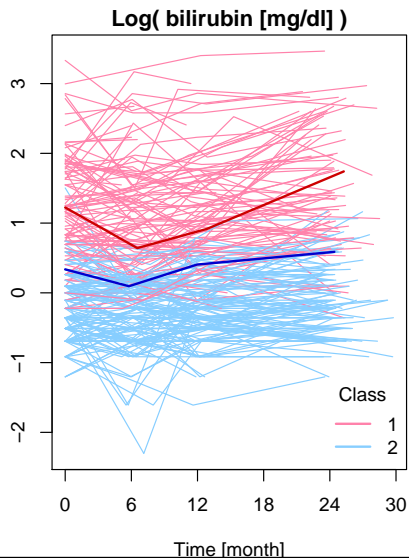
Subset PBC910 (see `mixAK` in )

Patient ID	Month	Gender	Drug	Log bilirubin	Spiders	Edema	...
3	0.00	male	D-penicillamine	0.34	0	0.5	...
	5.78	male	D-penicillamine	0.10	1	0	...
	11.96	male	D-penicillamine	0.41	0	0.5	...
	24.41	male	D-penicillamine	0.59	1	0.5	...
5	0.00	female	placebo	1.22	1	0	...
	6.54	female	placebo	0.64	0	0	...
	12.85	female	placebo	0.92	0	0.5	...
	25.26	female	placebo	1.74	1	1	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	

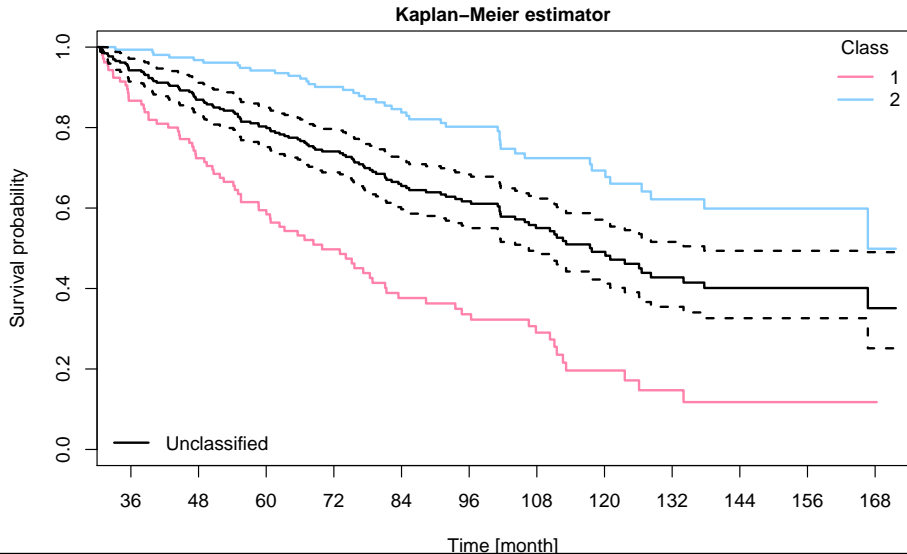
Mayo Clinic PBC Data - longitudinal profiles



Mayo Clinic PBC Data - after classification



Mayo Clinic PBC Data - usage of classification



Model based clustering

- Origins: Banfield and Raftery (1993)
- Outcomes: $\mathbf{Y}_i, i = 1, \dots, n$
- K models: $f_k(\mathbf{y}_i; \mathbf{x}_i, \psi, \psi^{(k)}), k \in \{1, \dots, K\}$
- Latent group indicators: $U_i \in \{1, \dots, K\}, \mathbf{Y}_i | U_i = k \sim f_k$
- Group probabilities: $0 < w_k = \mathbf{P}[U_i = k]$
- Parameters of interest: $\theta = \{\mathbf{w}, \psi, \psi^{(1)}, \dots, \psi^{(K)}\}$
- Mixture likelihood:

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


Model based clustering - classification rule

- By Bayes Theorem:

$$p_{i,k}(\theta) = P[U_i = k | \mathbf{Y}_i = \mathbf{y}_i; \mathbf{x}_i, \theta] = \frac{w_k f_k(\mathbf{y}_i; \mathbf{x}_i, \psi, \psi^{(k)})}{\sum_{\ell=1}^K w_\ell f_\ell(\mathbf{y}_i; \mathbf{x}_i, \psi, \psi^{(\ell)})}$$

- Classification rule:

$$\hat{U}_i := k \iff k = \arg \max_{\ell \in \{1, \dots, K\}} \hat{p}_{i,\ell}$$

- Estimation: Bayesian approach and MCMC methods (Gibbs sampling)
- Software: implemented in  using the C programming language

Numeric outcome

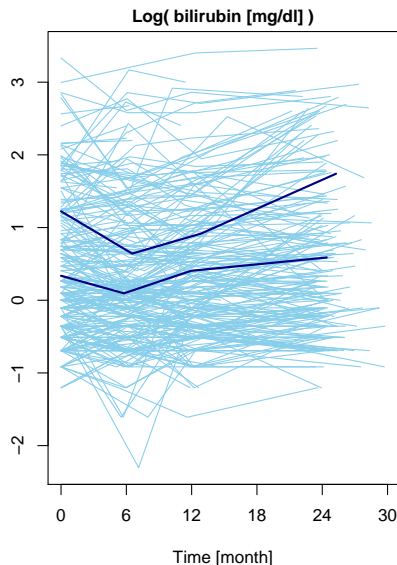
Logarithm of serum bilirubin [mg/dl]

Used model:

- Random effects models (LMM)
 - Laird and Ware (1983)
 - $N(\mathbf{X}_{ij}^T \boldsymbol{\beta} + \mathbf{Z}_{ij}^T \mathbf{b}_i, \sigma^2)$
 - $\mathbf{b}_i \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

For example:

- Fixed: $\beta_0 + \beta_1 \cdot \text{Month} + \dots$
 $\dots + \text{other regressors}$
- Random: $b_0 + b_1 \cdot \text{Month}$



Binary + Ordinal outcome

Presence and status of edema

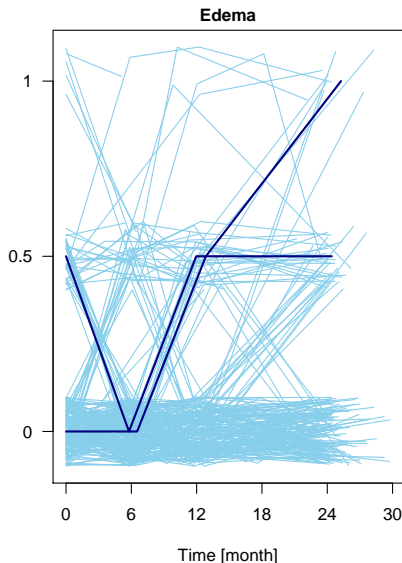
$L = 3$ ordered categories:

- 1 = Edema despite diuretic therapy,
- 0.5 = Untreated/Succesfully treated,
- 0 = No edema.

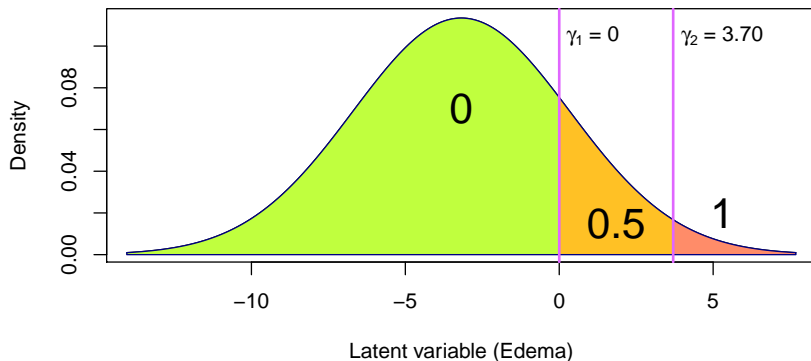
Used model

- **Latent variable modelling:** $Y|Y^*$
 - Y^* latent numeric outcome
 - Thresholding by

$$-\infty = \gamma_0 < \gamma_1 < \dots < \gamma_{L-1} < \gamma_L = \infty$$



Binary + Ordinal outcome - latent variable modelling



- **Fixed** threshold: $\gamma_1 = 0$
 - **Estimate** other thresholds: γ_2, \dots
 - $Y_{ij}^* | \mathbf{X}_{ij}, \mathbf{Z}_{ij}; \mathbf{b}_i \sim N(\mathbf{X}_{ij}^\top \boldsymbol{\beta} + \mathbf{Z}_{ij}^\top \mathbf{b}_i, 1)$
- | | | |
|----------------|--------|-------------------------------------|
| $Y_{ij} = 0$ | \iff | $Y_{ij}^* \leq \gamma_1$ |
| $Y_{ij} = 0.5$ | \iff | $\gamma_1 < Y_{ij}^* \leq \gamma_2$ |
| $Y_{ij} = 1$ | \iff | $\gamma_2 < Y_{ij}^*$ |

Joint modelling

Models for individual outcomes

- Numeric: $Y_{ij}^N \mid \mathbf{X}_{ij}^N, \mathbf{Z}_{ij}^N; \mathbf{b}_i^N \sim N \left(\left(\mathbf{X}_{ij}^N \right)^\top \boldsymbol{\beta}^N + \left(\mathbf{Z}_{ij}^N \right)^\top \mathbf{b}_i^N, (\tau^N)^{-1} \right)$
- Binary/Ordinal: $Y_{ij}^{O,*} \mid \mathbf{X}_{ij}^O, \mathbf{Z}_{ij}^O; \mathbf{b}_i^O \sim N \left(\left(\mathbf{X}_{ij}^O \right)^\top \boldsymbol{\beta}^O + \left(\mathbf{Z}_{ij}^O \right)^\top \mathbf{b}_i^O, 1 \right)$
 - $Y_{ij}^{O,*} \xrightarrow{\text{thresholding}} Y_{ij}^O$

Joining individual models through **joint distribution** of random effects:

$$\mathbf{b}_i = \begin{pmatrix} \mathbf{b}_i^N \\ \mathbf{b}_i^O \end{pmatrix} \sim N \left(\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}^N \\ \boldsymbol{\mu}^O \end{pmatrix}, \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}^{NN} & \boldsymbol{\Sigma}^{NO} \\ \boldsymbol{\Sigma}^{ON} & \boldsymbol{\Sigma}^{OO} \end{pmatrix} \right) \quad \text{leading to}$$

$$\text{var} \left[\begin{pmatrix} Y_{i,j}^N \\ Y_{i,j}^{O,*} \end{pmatrix} \mid C_{i,j} \right] = \begin{pmatrix} (\tau^N)^{-1} + \left(\mathbf{Z}_{i,j}^N \right)^\top \boldsymbol{\Sigma}^{NN} \mathbf{Z}_{i,j}^N & \left(\mathbf{Z}_{i,j}^N \right)^\top \boldsymbol{\Sigma}^{NO} \mathbf{Z}_{i,j}^O \\ \left(\mathbf{Z}_{i,j}^O \right)^\top \boldsymbol{\Sigma}^{ON} \mathbf{Z}_{i,j}^N & 1 + \left(\mathbf{Z}_{i,j}^O \right)^\top \boldsymbol{\Sigma}^{OO} \mathbf{Z}_{i,j}^O \end{pmatrix}$$

Final model

$$\mathbf{y}_i = \begin{pmatrix} \mathbf{y}_i^{\text{N}} \\ \mathbf{y}_i^{\text{O}} \end{pmatrix}, \quad \mathbf{y}_i^* = \begin{pmatrix} \mathbf{y}_i^{\text{N}} \\ \mathbf{y}_i^{\text{O},*} \end{pmatrix}, \quad \mathcal{C}_i = \left\{ \begin{array}{cccccc} \mathbf{x}_{i1}^{\text{N}} & \cdots & \mathbf{x}_{in_i}^{\text{N}} & \mathbf{z}_{i1}^{\text{N}} & \cdots & \mathbf{z}_{in_i}^{\text{N}} \\ \mathbf{x}_{i1}^{\text{O}} & \cdots & \mathbf{x}_{in_i}^{\text{O}} & \mathbf{z}_{i1}^{\text{O}} & \cdots & \mathbf{z}_{in_i}^{\text{O}} \end{array} \right\}$$

Choose class-specific parameters: $\psi^{(k)} = \{\beta^{(k)}, \mu^{(k)}, \Sigma^{(k)}\}$

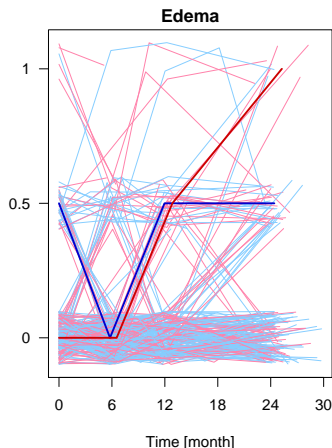
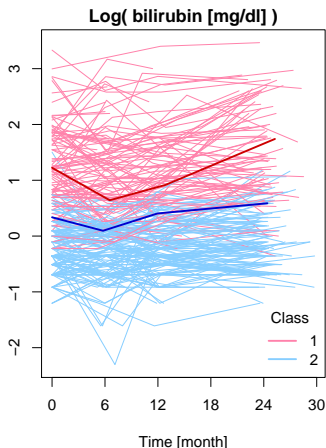
PDF of k -th model:

$$f_k(\mathbf{y}_i | \mathcal{C}_i; \gamma, \tau, \beta^{(k)}, \mu^{(k)}, \Sigma^{(k)}) = \int \int \rho(\mathbf{y}_i^{\text{O}} | \mathbf{y}_i^{\text{O},*}; \gamma) \cdot \\ \cdot \rho(\mathbf{y}_i^* | \mathcal{C}_i, \mathbf{b}_i; \beta^{(k)}, \tau) \cdot \rho(\mathbf{b}_i | \mu^{(k)}, \Sigma^{(k)}) d\mathbf{b}_i d\mathbf{y}_i^{\text{O},*}$$

MLE complicated \longrightarrow MCMC estimation

Classification probability

$$p_{i,k}(\theta) = P[U_i = k | \mathbf{Y}_i = \mathbf{y}_i; \mathbf{x}_i, \theta] = \frac{w_k f_k(\mathbf{y}_i; \mathbf{x}_i, \psi, \psi^{(k)})}{\sum_{\ell=1}^K w_\ell f_\ell(\mathbf{y}_i; \mathbf{x}_i, \psi, \psi^{(\ell)})}$$



Classification rule of new patient

New patient: $\mathbf{Y}_{\text{new}}, \mathbf{X}_{\text{new}}, \mathbf{Z}_{\text{new}}$ and unobserved $U_{\text{new}}, \mathbf{b}_{\text{new}}, \mathbf{Y}_{\text{new}}^{\text{O},\star}$

Bayesian estimator: $\hat{p}_{\text{new},k} = E [p_{\text{new},k}(\boldsymbol{\theta}) \mid \text{data } \mathbb{Y}, \mathcal{C}]$

Calculation: use generated sample $\boldsymbol{\theta}^1, \dots, \boldsymbol{\theta}^M$ from $p(\boldsymbol{\theta} \mid \text{data } \mathbb{Y}, \mathcal{C})$

Gradual classification:

$$\begin{array}{rcccl}
 \mathbf{y}_{\text{new},1} & & & \longrightarrow & p_{\text{new},k}^1 \\
 \mathbf{y}_{\text{new},1} & \mathbf{y}_{\text{new},2} & & \longrightarrow & p_{\text{new},k}^2 \\
 \vdots & \vdots & \ddots & & \vdots \\
 \mathbf{y}_{\text{new},1} & \mathbf{y}_{\text{new},2} & \cdots & \mathbf{y}_{\text{new},n_{\text{new}}} & \longrightarrow & p_{\text{new},k}^{n_{\text{new}}}
 \end{array}$$

Further development

- GLMM as an alternative to latent variable modelling.
- Add general categorical variables.
- Perfecting MCMC algorithm.
- Robustness to violation of normality assumption.
- Selection of important regressors.

Further development

- GLMM as an alternative to latent variable modelling.
- Add general categorical variables.
- Perfecting MCMC algorithm.
- Robustness to violation of normality assumption.
- Selection of important regressors.

Any interesting suggestions?

Further development

- GLMM as an alternative to latent variable modelling.
- Add general categorical variables.
- Perfecting MCMC algorithm.
- Robustness to violation of normality assumption.
- Selection of important regressors.

Any interesting suggestions?

Thank you for your attention.

Literature references

- [1] D. Banfield, J. and E. Raftery, A. Model-based Gaussian and non-Gaussian clustering. *Biometrics*, 49(3):803–821, 1993.
- [2] R. Fleming, T. and P. Harrington, D. *Counting Processes and Survival Analysis*. New York: John Wiley and Sons, New York, 1991.
- [3] A. Komárek and L. Komárková. Capabilities of R package mixAK for clustering based on multivariate continuous and discrete longitudinal data. *Journal of Statistical Software*, 59(12):1–38, 2014.
- [4] M. Laird, N. and H. Ware, J. Random-effects models for longitudinal data. *Biometrics*, 38:963–74, 01 1983.
- [5] C. P. Robert and G. Casella. *Monte Carlo statistical methods*. Second edition. Springer-Verlag, New York, 2004.
- [6] S. Wilhelm and B. Manjunath. tmvtnorm: A package for the truncated multivariate normal distribution. *The R Journal*, 2, 06 2010.