

# Sdružené modelování spojитých i diskrétních longitudinálních dat s exkurzí do diskriminační a shlukové analýzy

Arnošt Komárek

Katedra pravděpodobnosti a matematické statistiky

Matematicko-fyzikální fakulta  
Univerzity Karlovy v Praze

XVI. ROBUST 2010

Králíky, 31. ledna – 5. února 2010

- ① Zpět do Nečtin a taky trochu do Roháčů
- ② Na skok do Hejnic
- ③ Zpátky do Králík
- ④ Stále v Králíkách
- ⑤ Co v Králíkách bude možná až po 22. hodině + další dvouletka
- ⑥ Něco na závěr

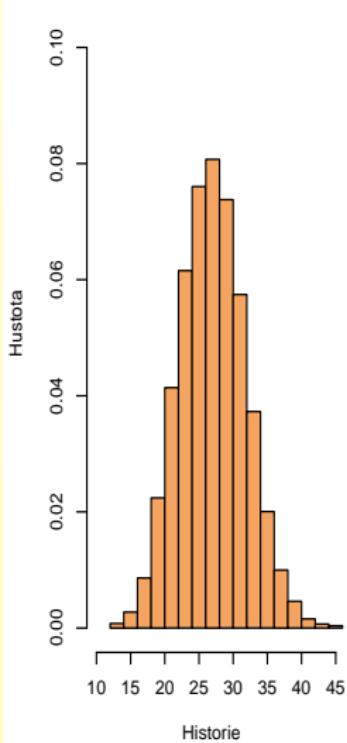
## Část I

Zpět do Nečtin a taky trochu do Roháčů

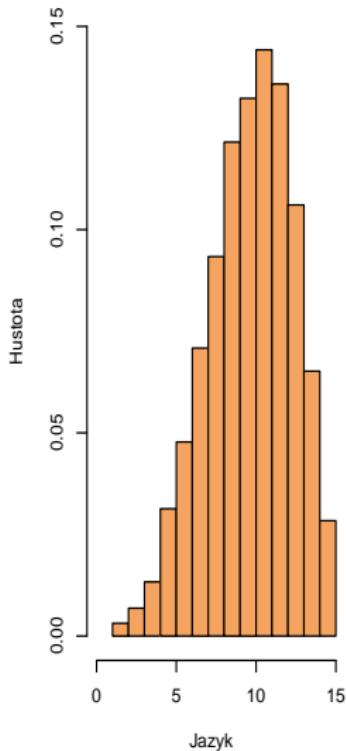
- ❖ *Matematika mezi §§§ aneb něco málo o diskriminaci*
- ❖ 12. termín přijímaček na Právnickou fakultu UK v Praze v roce 1999 s únikem (a prodejem) zadání testů  
*(které byly dle tehdejšího rektora UK i děkana PF dílem gangsterské mafie stojící mimo fakultu)*
- ❖ Lze vtipovat uchazeče napojené na tuto mafii?

# Přijímačky na PF UK v roce 1999 (Termíny č. 1–11)

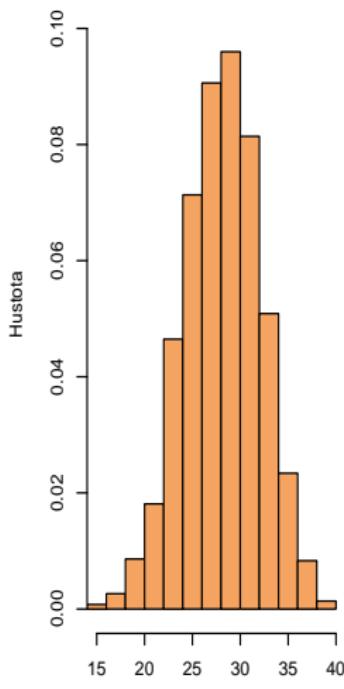
Historie



Jazyk

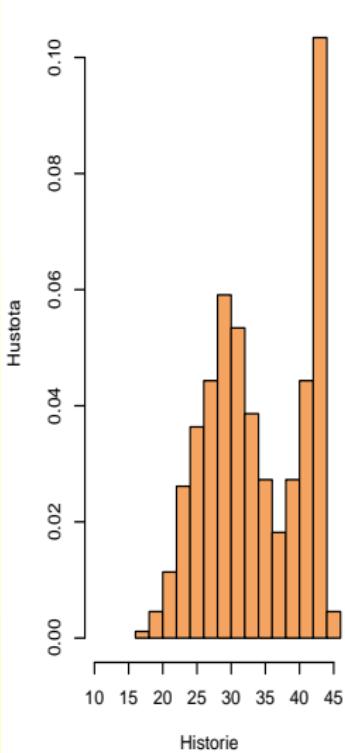


Logika

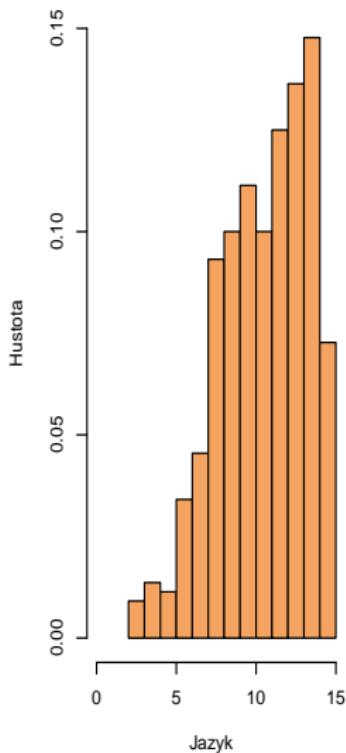


# Přijímačky na PF UK v roce 1999 (Termín č. 12)

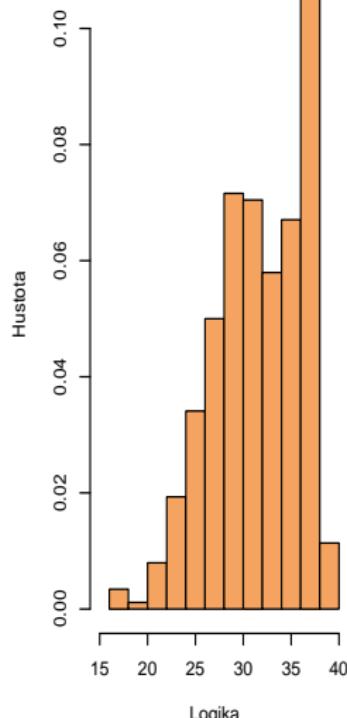
Historie



Jazyk

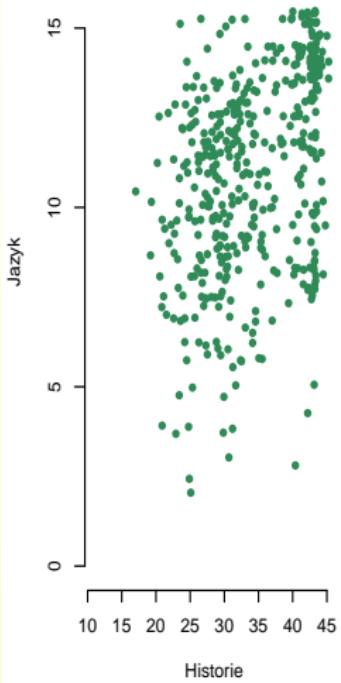


Logika

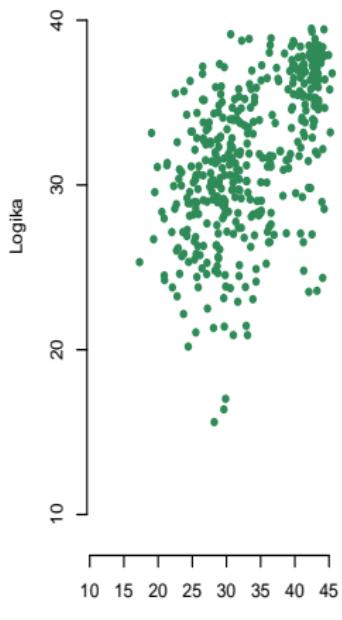


# Přijímačky na PF UK v roce 1999 (Termín č. 12)

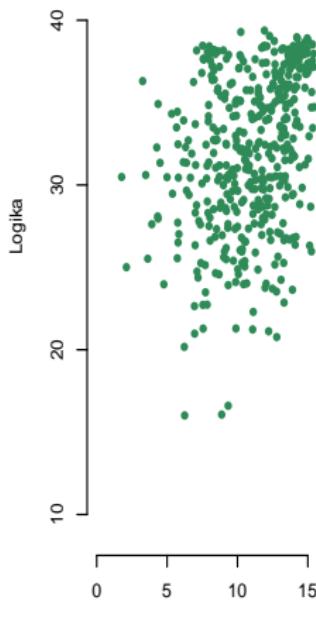
Historie x Jazyk



Historie x Logika



Jazyk x Logika



# Modelově založená shluková analýza (Model based clustering)

- ❖  $\mathbf{Y} = (Y_1, \dots, Y_p)'$  : náhodný vektor pro odezvu  
(např. body z jednotlivých testů u náhodně vybraného uchazeče) s hustotou  $p(\mathbf{y})$
- ❖ Směsový model pro  $\mathbf{Y}$  :  $p(\mathbf{y}) = \sum_{k=1}^K w_k \phi(\mathbf{y}; \boldsymbol{\eta}_k)$
- ❖  $K$  ... počet shluků
- ❖  $\mathbf{w} = (w_1, \dots, w_K)'$  ... proporce (váhy) jednotlivých shluků
- ❖  $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_K$  ... parametry rozdělení odezvy v jednotlivých shlucích
- ❖  $\phi$  ... vhodná hustota
- ❖  $\boldsymbol{\theta} = (w_1, \dots, w_K, \boldsymbol{\eta}'_1, \dots, \boldsymbol{\eta}'_K)'$  ... parametry modelu

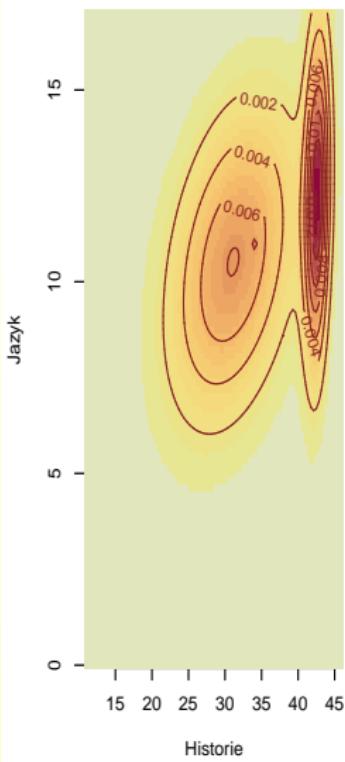
# Modelově založená shluková analýza

Odhad směsového modelu

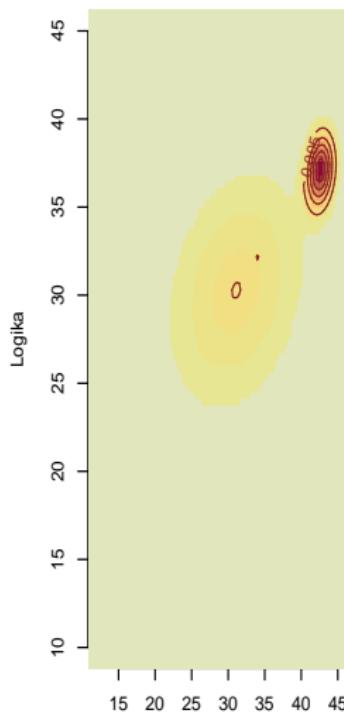
- ❖  $\mathbf{Y}_1, \dots, \mathbf{Y}_n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{y}) = \sum_{k=1}^K w_k \phi(\mathbf{y}; \boldsymbol{\eta}_k)$
- ❖ Nečtiny (XI. ROBUST):  $\phi(\mathbf{y}; \boldsymbol{\eta}_k) \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma})$ 
  - ✿ Odhad metodou maximální věrohodnosti (EM algoritmus)
- ❖ Roháče (XV. ROBUST):  $\phi(\mathbf{y}; \boldsymbol{\eta}_k) \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ 
  - ✿ Bayesovský odhad přes MCMC

# Přijímačky na PF UK v roce 1999 (Termín č. 12)

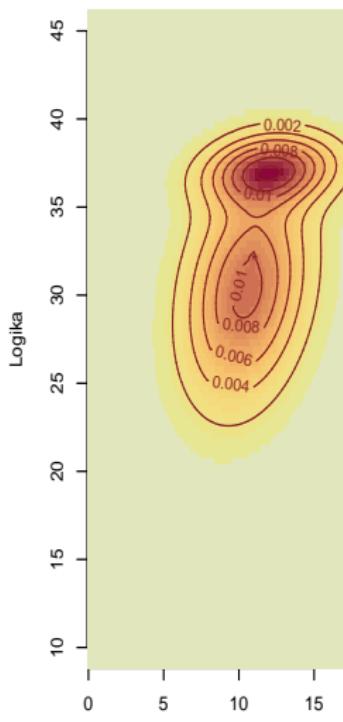
Historie x Jazyk



Historie x Logika



Jazyk x Logika



# Modelově založená shluková analýza

## Indikátory shluků

- ❖ Pro shlukování: i.i.d. náhodné veličiny  $U_1, \dots, U_n$
- ❖  $U_i \in \{1, \dots, K\}$
- ❖  $P(U_i = k) = w_k, k = 1, \dots, K, i = 1, \dots, n$

# Modelově založená shluková analýza

Směsový model zapsaný hierarchicky

❖  $\mathbf{Y}_i \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{y}) = \sum_{k=1}^K w_k \phi(\mathbf{y}; \boldsymbol{\eta}_k)$

❖ Hierarchicky:

$$\left. \begin{array}{rcl} P(U_i = k) & = & w_k \\ p(\mathbf{y}_i | U_i = k) & = & \phi(\mathbf{y}_i; \boldsymbol{\eta}_k) \end{array} \right\} \quad k = 1, \dots, K$$

---

► Bayesova věta:

$$\pi_{i,k}(\boldsymbol{\theta}) \equiv P(U_i = k | \mathbf{y}_i) = \frac{w_k \phi(\mathbf{y}_i; \boldsymbol{\eta}_k)}{\sum_{j=1}^K w_j \phi(\mathbf{y}_i; \boldsymbol{\eta}_j)}$$

# Modelově založená shluková analýza

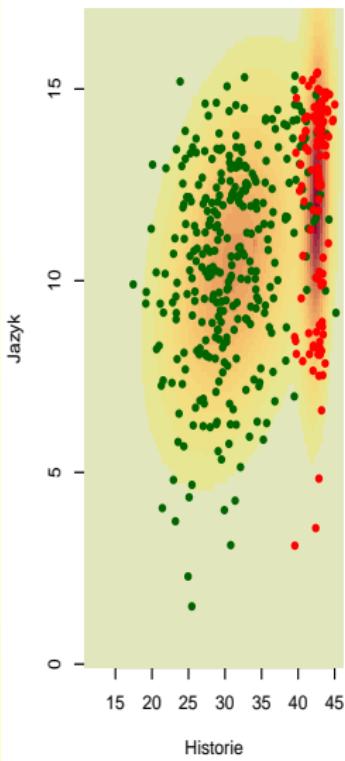
## Shlukování

- Zařadí  $i$ -té pozorování do shluku s indexem  $j$ :

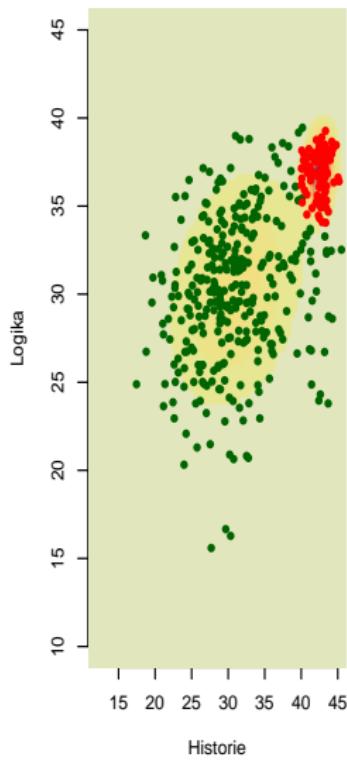
$$\pi_{i,j} = \max\{\pi_{i,1}, \dots, \pi_{i,K}\}$$

# Přijímačky na PF UK v roce 1999 (Termín č. 12)

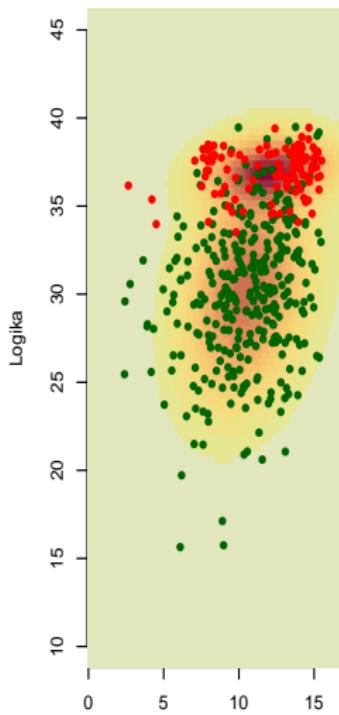
Historie x Jazyk



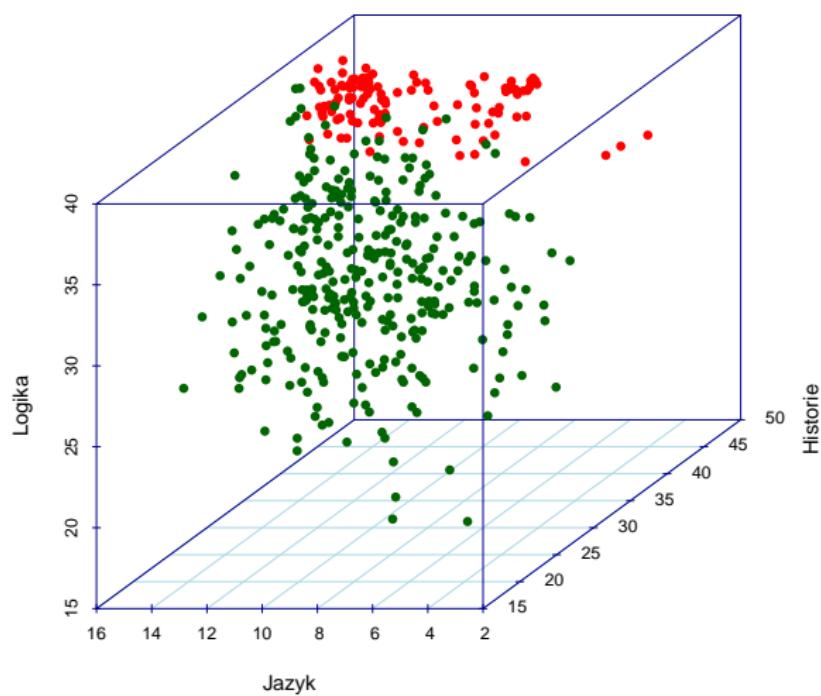
Historie x Logika



Jazyk x Logika



# Přijímačky na PF UK v roce 1999 (Termín č. 12)



## Část II

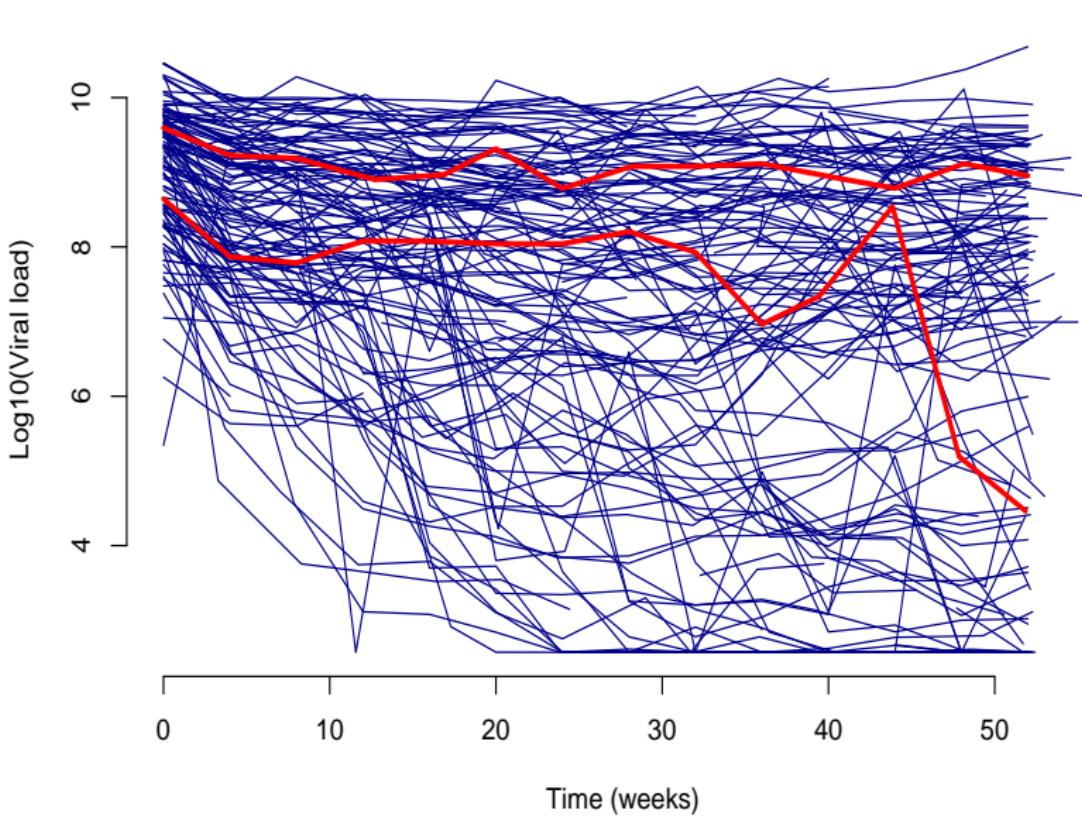
Na skok do Hejnic

# Chronická hepatitida B

- ❖ Mezinárodní klinická studie koordinovaná na Erasmus Medisch Centrum v Rotterdamu
- ❖ Léčba pomocí PEG-INF (na ROBUSTu nekombinováno s lamivudinem) po dobu 52 týdnů
- ❖ Má to nehezké vedlejší účinky a často je stejně k ničemu
- ❖ Podrobnosti:  
Janssen, van Zonneveld, Senturk, Akarca, Cakaloglu, Simon, So, Gerken, de Man, Niesters, Zondervan, Hansen, Schalm (2005).  
*Pegylated interferon alfa-2b alone or in combination with lamivudine for HBeAg-positive chronic hepatitis B: a randomised trial.* *Lancet*, **365**, 123–129.

# Chronická hepatitida B

Log10(Viral load)



# Chronická hepatitida B

- ❖ Snaha roztrídit pacienty do skupin
    - ❖ může nějak souviset s diagnózou či prognózou
    - ❖ ...
- 

- ❖ Jeden pacient:  $\mathbf{Y}_i = (Y_{i,1}, \dots, Y_{i,n_i})'$
  - ❖  $\mathbf{Y}_i \sim p(\mathbf{y})$
- 

- ❖ Liší se to nějak od právnické mafie?

# Model pro shlukování?

- ❖ Jeden pacient:  $\mathbf{Y}_i = (Y_{i,1}, \dots, Y_{i,n_i})'$
  - ❖  $n_i$  (počet návštěv) není nutně stejné pro všechny pacienty
- 
- ❖ Ve skutečnosti je  $Y_{i,j} = Y_{i,j}(t_{i,j})$ 
    - ❖  $t_{i,j}$  ... doba, po kterou je pacient léčen
  - ❖ Posloupnost časů návštěv  $(t_{i,1}, \dots, t_{i,n_i})$  není stejná pro všechny pacienty
- 
- ➡ Těžko bude  $\mathbf{Y}_i \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{y})$

# Model pro shlukování?

- ❖ Co takhle reprezentovat každého pacienta nějakým (nepřímo pozorovatelným) náhodným vektorem  $\mathbf{b}_i$ , tak, aby  $\mathbf{b}_i \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{b})$

- ❖ *On a fitting of a linear mixed model with a finite normal mixture as random-effects distribution*
- ❖ Lineární smíšený model (LMM) pro  $\mathbf{Y}_i$ :

$$\mathbf{Y}_i = \mathbb{X}_i \boldsymbol{\alpha} + \mathbb{Z}_i \mathbf{b}_i + \boldsymbol{\varepsilon}_i$$

- \*  $\mathbb{X}_i, \mathbb{Z}_i$  ... matice regresorů
- \*  $\boldsymbol{\alpha}$  ... (neznámé) regresní parametry
- \*  $\mathbf{b}_i \stackrel{\text{i.i.d.}}{\sim} \sum_{k=1}^K w_k \mathcal{N}(\boldsymbol{\mu}_k, \mathbb{D}_k)$
- \*  $\boldsymbol{\varepsilon}_i \sim \mathcal{N}(\mathbf{0}, \sigma^2 I_{n_i})$ , nezávislé pro  $i = 1, \dots, n$

❖ Hejnice (XII. ROBUST):

- ✳ Pouze homoskedastická směs
- ✳ Maximální věrohodnost přes EM algoritmus s využitím software pro standardní (nesměsový) LMM

❖ Králíky (XVI. ROBUST):

- ✳ Směs může být i heteroskedastická
- ✳ Bayesovský odhad přes MCMC
- ✳ Nestandardní podsoftware (balíček) standardního software (R)

# Modelově založená shluková analýza

## Indikátory shluků

- ❖ Pro shlukování: i.i.d. náhodné veličiny  $U_1, \dots, U_n$
- ❖  $U_i \in \{1, \dots, K\}$
- ❖  $P(U_i = k) = w_k, k = 1, \dots, K, i = 1, \dots, n$

# Modelově založená shluková analýza

LMM s normální směsí v rozdělení náhodných efektů zapsaný hierarchicky

❖ Hierarchicky zapsáno:

$$[\mathbf{Y}_i | \mathbf{b}_i] \sim \mathcal{N}(\mathbb{X}_i \boldsymbol{\alpha} + \mathbb{Z}_i \mathbf{b}_i, \sigma^2 I_{n_i})$$

$$\left. \begin{array}{rcl} \mathsf{P}(U_i = k) & = & w_k \\ [\mathbf{b}_i | U_i = k] & \sim & \mathcal{N}(\boldsymbol{\mu}_k, \mathbb{D}_k) \end{array} \right\} \quad k = 1, \dots, K$$

---

$$\Rightarrow \mathbf{Y}_i \sim \sum_{k=1}^K w_k \mathcal{N}(\mathbb{X}_i \boldsymbol{\alpha} + \mathbb{Z}_i \boldsymbol{\mu}_k, \mathbb{Z}_i \mathbb{D}_k \mathbb{Z}_i' + \sigma^2 I_{n_i})$$

$$\Rightarrow [\mathbf{Y}_i | U_i = k] \sim \mathcal{N}(\mathbb{X}_i \boldsymbol{\alpha} + \mathbb{Z}_i \boldsymbol{\mu}_k, \mathbb{Z}_i \mathbb{D}_k \mathbb{Z}_i' + \sigma^2 I_{n_i})$$

# Chronická hepatitida B

- ❖  $Y_{i,j} = b_{i,1} + b_{i,2}B_2(t_{i,j}) + b_{i,3}B_3(t_{i,j}) + \varepsilon_{i,j}$
- ❖  $1, B_2(t), B_3(t)$   
... kvadratická B-splinová báze s uzly v  $t = 0$  a  $56$

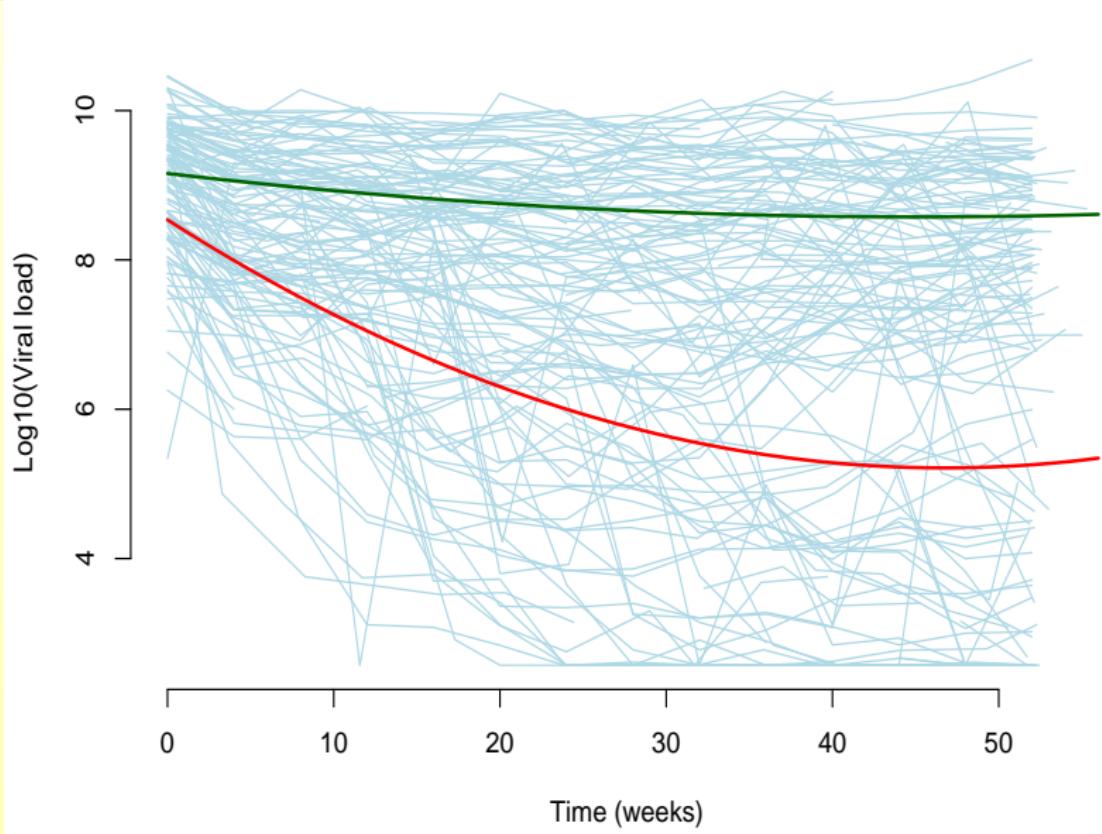
---

- ❖  $\boldsymbol{b}_i = (b_{i,1}, b_{i,2}, b_{i,3})'$  <sup>i.i.d.</sup>  $\sim \sum_{k=1}^2 w_k \mathcal{N}(\boldsymbol{\mu}_k, \mathbb{D}_k)$

⇒  $E(Y_{i,j} \mid U_i = k) = \mu_{k,1} + \mu_{k,2}B_2(t_{i,j}) + \mu_{k,3}B_3(t_{i,j})$

# Chronická hepatitida B

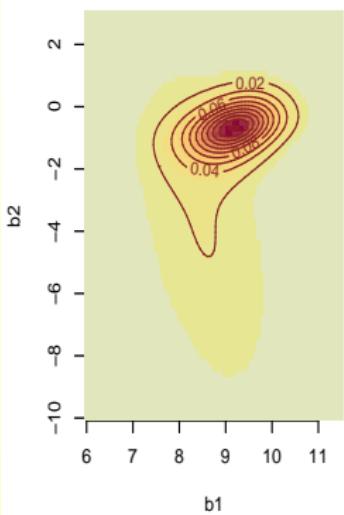
Odhadnutý průměrný vývoj v čase ve dvou skupinách



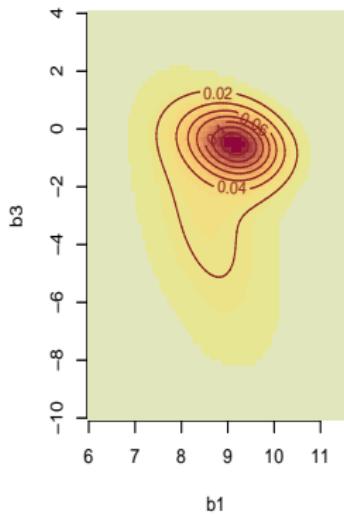
# Chronická hepatitida B

Rozdělení náhodných efektů

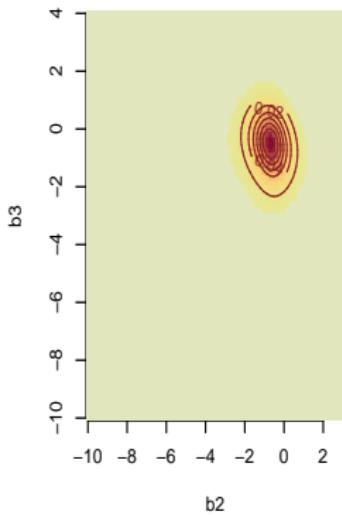
$b_1 \times b_2$



$b_1 \times b_3$



$b_2 \times b_3$



# Modelově založená shluková analýza

## ❖ Model:

$$[\mathbf{Y}_i | \mathbf{b}_i] \sim \mathcal{N}(\mathbb{X}_i \boldsymbol{\alpha} + \mathbb{Z}_i \mathbf{b}_i, \sigma^2 I_{n_i})$$

$$\left. \begin{array}{lcl} \mathsf{P}(U_i = k) & = & w_k \\ [\mathbf{b}_i | U_i = k] & \sim & \mathcal{N}(\boldsymbol{\mu}_k, \mathbb{D}_k) \end{array} \right\} \quad k = 1, \dots, K$$

## ❖ Parametry modelu:

$$\boldsymbol{\theta} = (\boldsymbol{\alpha}', \mathbf{w}', \boldsymbol{\mu}'_1, \dots, \boldsymbol{\mu}'_K, \text{vec}(\mathbb{D}_1)', \dots, \text{vec}(\mathbb{D}_K)', \sigma^2)' \in \Theta$$

## ❖ Nepřímo pozorovatelné náhodné veličiny/vektory:

$$\boldsymbol{\eta} = (\mathbf{b}'_1, \dots, \mathbf{b}'_n, \mathbf{U}_1, \dots, \mathbf{U}_n)' \in \mathbb{R}^q \times \dots \times \mathbb{R}^q \times \{1, \dots, K\}^n$$

# Modelově založená shluková analýza

→ Bayesova věta:

$$\begin{aligned}\pi_{i,k}(\theta) &\equiv P(U_i = k \mid \mathbf{y}_i, \theta) \\ &\propto w_k p(\mathbf{y}_i \mid U_i = k, \theta) \\ &= w_k \varphi(\mathbf{y}_i; \mathbb{X}_i \boldsymbol{\alpha} + \mathbb{Z}_i \boldsymbol{\mu}_k, \mathbb{Z}_i \mathbb{D}_k \mathbb{Z}_i' + \sigma^2 I_{n_i})\end{aligned}$$

→ Po chvíli počítání:

$$\begin{aligned}\pi_{i,k}(\theta) &\propto \\ w_k |\mathbb{D}_k|^{-1/2} |\mathbb{A}_{i,k}|^{-1/2} \exp\left\{-\frac{1}{2} (\boldsymbol{\mu}_k' \mathbb{D}_k^{-1} \boldsymbol{\mu}_k - \mathbf{c}'_{i,k} \mathbb{A}_{i,k}^{-1} \mathbf{c}_{i,k})\right\} \\ \mathbb{A}_{i,k} &= \sigma^{-2} \mathbb{Z}_i' \mathbb{Z}_i + \mathbb{D}_k^{-1} \\ \mathbf{c}_{i,k} &= \sigma^{-2} \mathbb{Z}_i' (\mathbf{y}_i - \mathbb{X}_i \boldsymbol{\alpha}) + \mathbb{D}_k^{-1} \boldsymbol{\mu}_k\end{aligned}$$

# Modelově založená shluková analýza

## Shlukování

- ❖ Je-li  $\hat{\theta}$  ML odhad  $\theta$  :
  - ✳ Shlukuj na základě  $\pi_{i,k}(\hat{\theta})$
- ❖ Je-li při Bayesovském odhadu,  $\theta^{(1)}, \dots, \theta^{(M)}$  náhodný výběr z aposteriorního rozdělení  $p(\theta | \mathbf{y}_1, \dots, \mathbf{y}_n)$  :
  - ✳ Shlukuj na základě

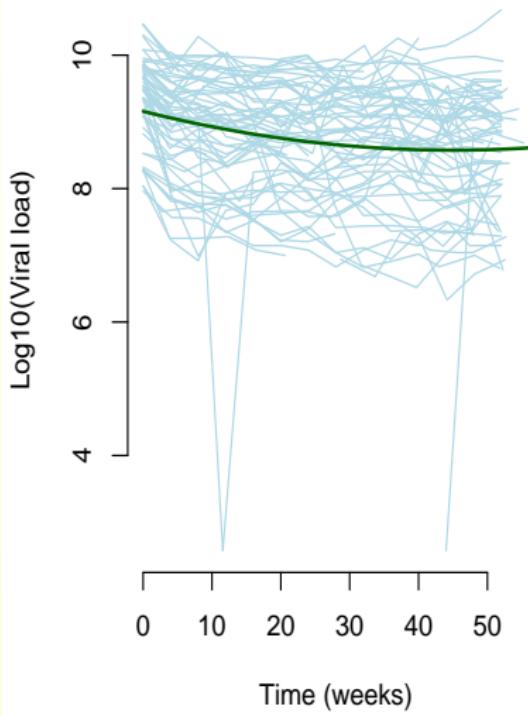
$$\widehat{P}(U_i = k | \mathbf{y}_1, \dots, \mathbf{y}_n) = \widehat{E}_{\theta}\{\pi_{i,k}(\theta) | \mathbf{y}_1, \dots, \mathbf{y}_n\}$$

$$= M^{-1} \sum_{m=1}^M \pi_{i,k}(\theta^{(m)})$$

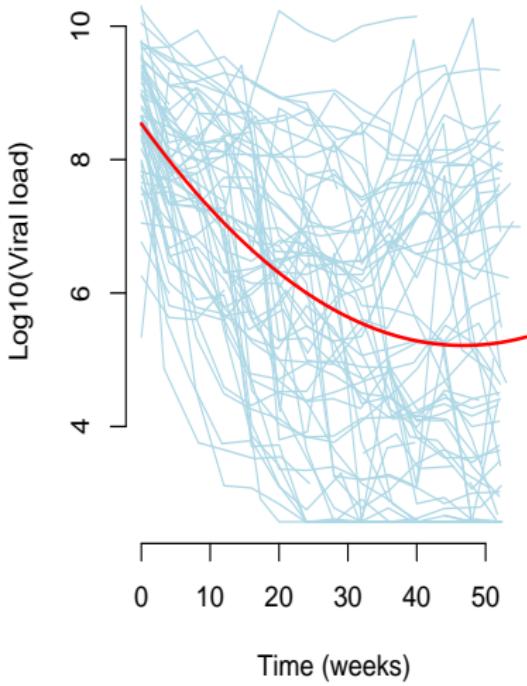
# Chronická hepatitida B

## Shluky dle modelu pro Log10(Viral load)

Model based group 0



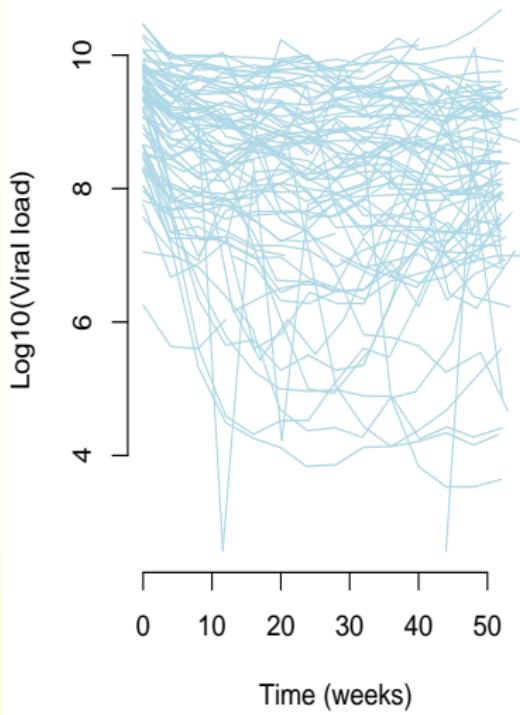
Model based group 1



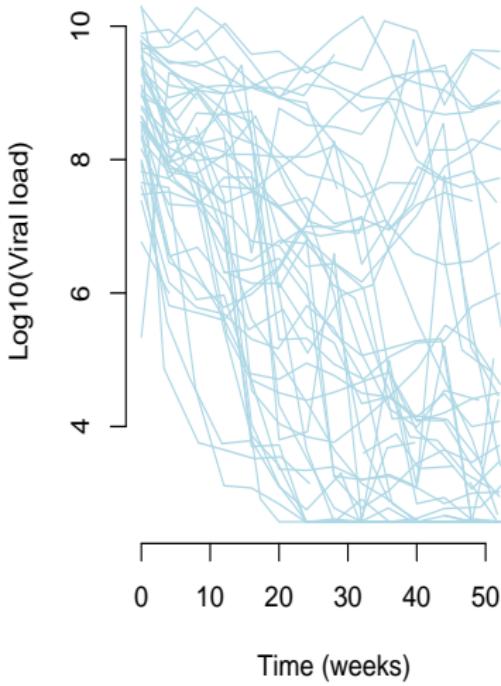
# Chronická hepatitida B

Shluky dle HBeAg v 78. týdnu

HBeAg78 group 0



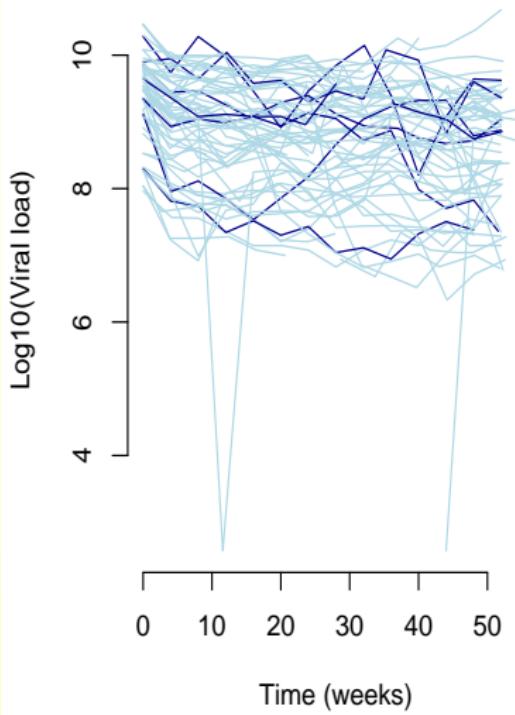
HBeAg78 group 1



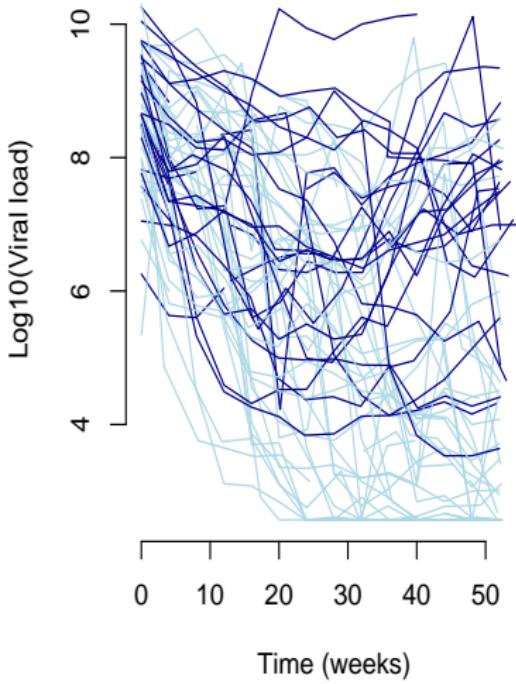
# Chronická hepatitida B

Shluky dle modelu pro Log10(Viral load), rozdíl oproti HBeAg78 shlukům

Model based group 0



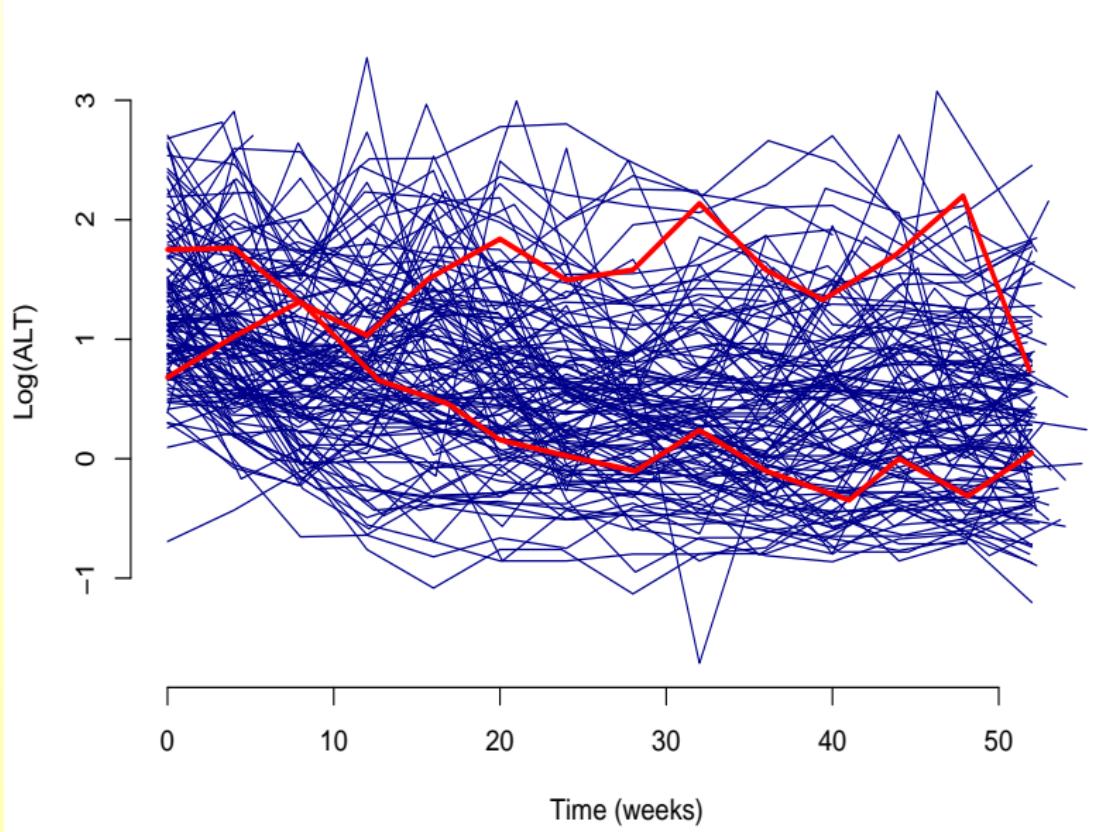
Model based group 1



- ❖ Nešlo by zlepšit odhalení skutečné HBeAg78 skupiny?
  - ❖ Kromě množství viru v krvi se měřilo ještě ALT  
( míra aktivity choroby)
- 
- ➡ Co takhle shlukovat na základě **obou** ukazatelů?

# Chronická hepatitida B

Log(ALT)



## Část III

### Zpátky do Králík

# LMM pro dvě různé odezvy

- ❖ Jeden pacient:

$$\mathbf{Y}_{i,1} = (Y_{i,1,1}, \dots, Y_{i,1,n_{i,1}})' \quad \text{Log10(Viral load)}$$

$$\mathbf{Y}_{i,2} = (Y_{i,2,1}, \dots, Y_{i,2,n_{i,2}})' \quad \text{Log(ALT)}$$

- ❖ LMM pro odezvu 1 (Log10(Viral load)):

$$\mathbf{Y}_{i,1} = \mathbb{X}_{i,1}\boldsymbol{\alpha}_1 + \mathbb{Z}_{i,1}\mathbf{b}_{i,1} + \boldsymbol{\varepsilon}_{i,1}$$

- ❖ LMM pro odezvu 2 (Log(ALT)):

$$\mathbf{Y}_{i,2} = \mathbb{X}_{i,2}\boldsymbol{\alpha}_2 + \mathbb{Z}_{i,2}\mathbf{b}_{i,2} + \boldsymbol{\varepsilon}_{i,2}$$

# LMM pro dvě různé odezvy

❖ Lineární smíšené modely:

$$\mathbf{Y}_{i,1} = \mathbb{X}_{i,1}\boldsymbol{\alpha}_1 + \mathbb{Z}_{i,1}\mathbf{b}_{i,1} + \boldsymbol{\varepsilon}_{i,1}$$

$$\mathbf{Y}_{i,2} = \mathbb{X}_{i,2}\boldsymbol{\alpha}_2 + \mathbb{Z}_{i,2}\mathbf{b}_{i,2} + \boldsymbol{\varepsilon}_{i,2}$$

- ❖  $\boldsymbol{\varepsilon}_{i,1} \sim \mathcal{N}(\mathbf{0}, \sigma_1^2 I_{n_{i,1}})$ , nezávislé pro  $i = 1, \dots, n$
- ❖  $\boldsymbol{\varepsilon}_{i,2} \sim \mathcal{N}(\mathbf{0}, \sigma_2^2 I_{n_{i,2}})$ , nezávislé pro  $i = 1, \dots, n$
- ❖  $\mathbf{b}_i = (b_{i,1,1}, \dots, b_{i,1,q_1}, b_{i,2,1}, \dots, b_{i,2,q_2})'$
- ❖  $\mathbf{b}_i \stackrel{\text{i.i.d.}}{\sim} \sum_{k=1}^K w_k \mathcal{N}(\boldsymbol{\mu}_k, \mathbb{D}_k)$

# LMM pro dvě různé odezvy napsaný jako jeden LMM

❖ Lze to napsat najednou:

$$\begin{pmatrix} \mathbf{Y}_{i,1} \\ \mathbf{Y}_{i,2} \end{pmatrix} = \begin{pmatrix} \mathbb{X}_{i,1} & \mathbb{O} \\ \mathbb{O} & \mathbb{X}_{i,2} \end{pmatrix} \begin{pmatrix} \boldsymbol{\alpha}_1 \\ \boldsymbol{\alpha}_2 \end{pmatrix} + \begin{pmatrix} \mathbb{Z}_{i,1} & \mathbb{O} \\ \mathbb{O} & \mathbb{Z}_{i,2} \end{pmatrix} \begin{pmatrix} \mathbf{b}_{i,1} \\ \mathbf{b}_{i,2} \end{pmatrix} + \begin{pmatrix} \boldsymbol{\varepsilon}_{i,1} \\ \boldsymbol{\varepsilon}_{i,2} \end{pmatrix}$$

❖  $\begin{pmatrix} \boldsymbol{\varepsilon}_{i,1} \\ \boldsymbol{\varepsilon}_{i,2} \end{pmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} \sigma_1^2 I_{n_{i,1}} & \mathbb{O} \\ \mathbb{O} & \sigma_2^2 I_{n_{i,2}} \end{pmatrix}\right)$ ,

nezávislé pro  $i = 1, \dots, n$

❖  $\mathbf{b}_i = (b_{i,1,1}, \dots, b_{i,1,q_1}, b_{i,2,1}, \dots, b_{i,2,q_2})'$

❖  $\mathbf{b}_i \stackrel{\text{i.i.d.}}{\sim} \sum_{k=1}^K w_k \mathcal{N}(\boldsymbol{\mu}_k, \mathbb{D}_k)$

# LMM pro dvě různé odezvy

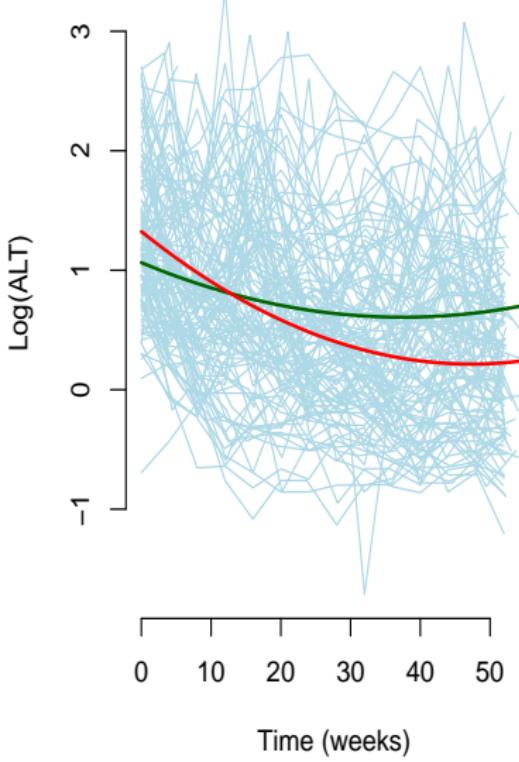
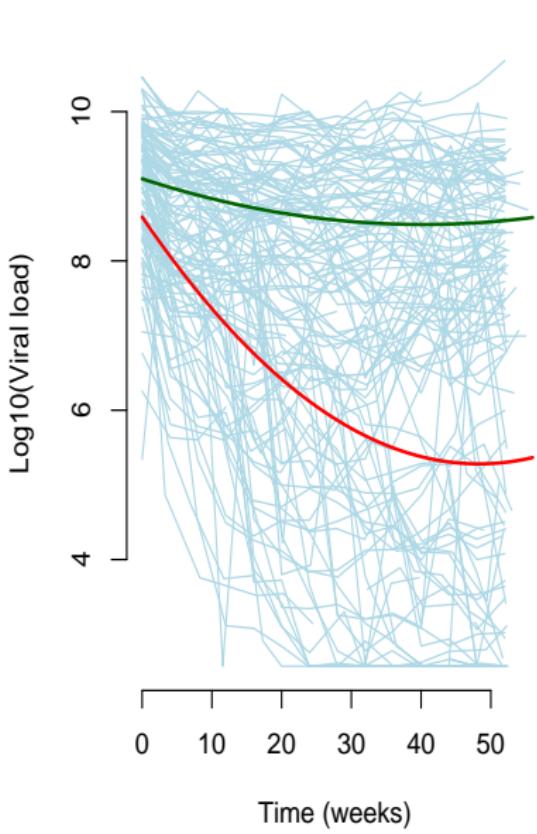
- Při  $K = 1$  lze ML odhadnout standardním softwarem  
(R balíčky nlme, lme4, SAS PROC MIXED)
  - Při  $K > 1$  a homoskedastické směsi v rozdělení náhodných efektů lze teoreticky použít přístup prezentovaný v Hejnicích na XII. ROBUSTu  
(EM algoritmus s M-krokem provedeným standardním softwarem)
- 

- ❖ Při použití standardního software se zbytečně počítá s mnoha nulami v maticích  $\mathbb{X}$ ,  $\mathbb{Z}$
- Dlouho to trvá
- Přináší numerické problémy (i při  $K = 1$ )

- ❖ Směs v rozdělení náhodných efektů může být i heteroskedastická
- ❖ Bayesovský odhad přes MCMC
- ❖ Nestandardní podsoftware (balíček) standardního software (R)
- ❖ Využití známé struktury matic  $\mathbb{X}$ ,  $\mathbb{Z}$

# Chronická hepatitida B

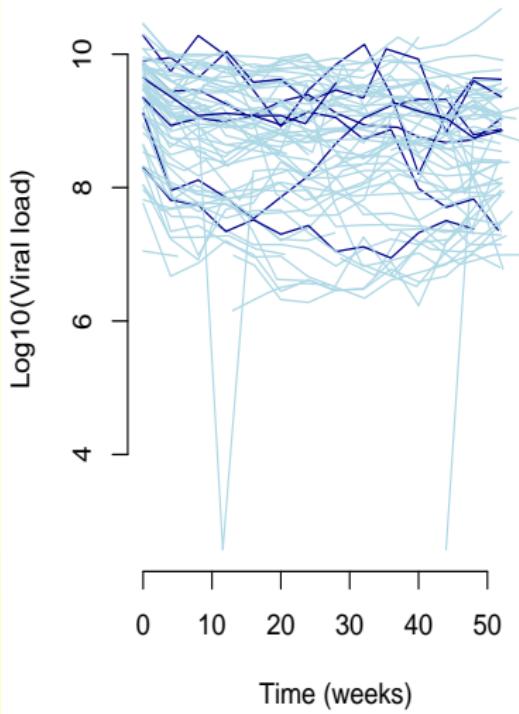
Odhadnutý průměrný vývoj v čase ve dvou skupinách



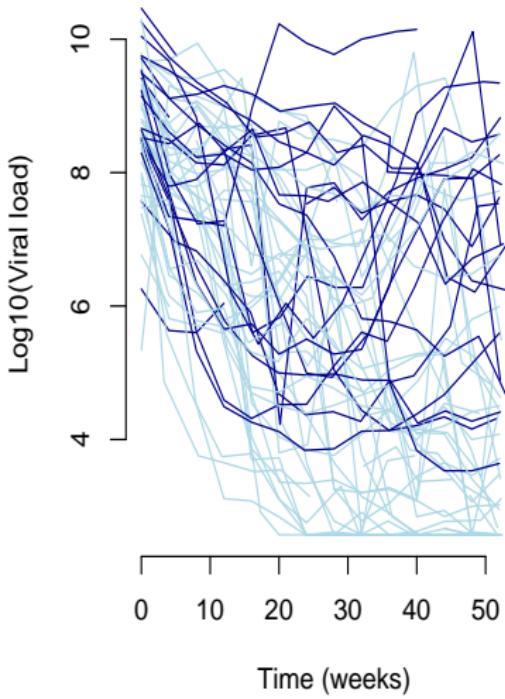
# Chronická hepatitida B

Shluky dle sdruženého modelu pro Log10(Viral load) a Log(ALT)

Model based group 0



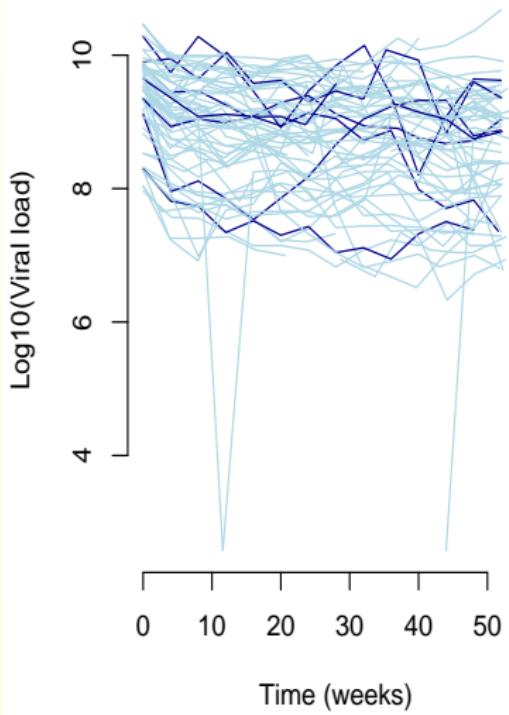
Model based group 1



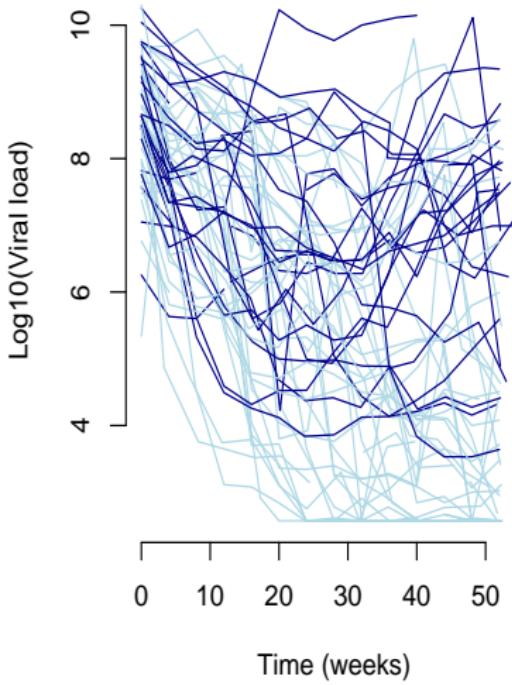
# Chronická hepatitida B

## Shluky dle modelu pro Log10(Viral load)

Model based group 0



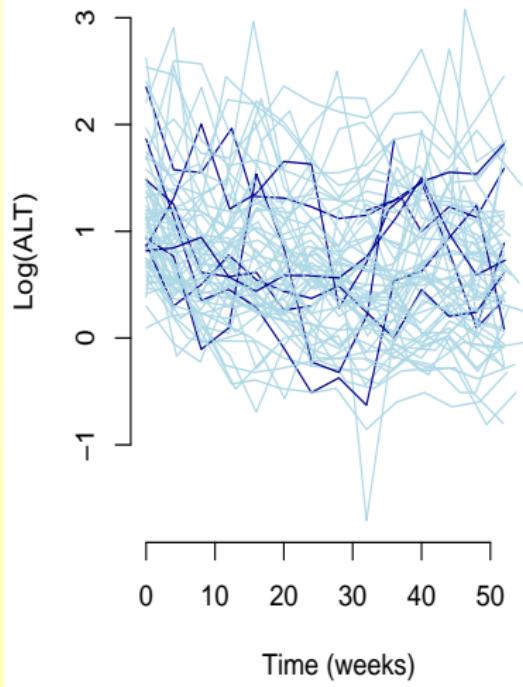
Model based group 1



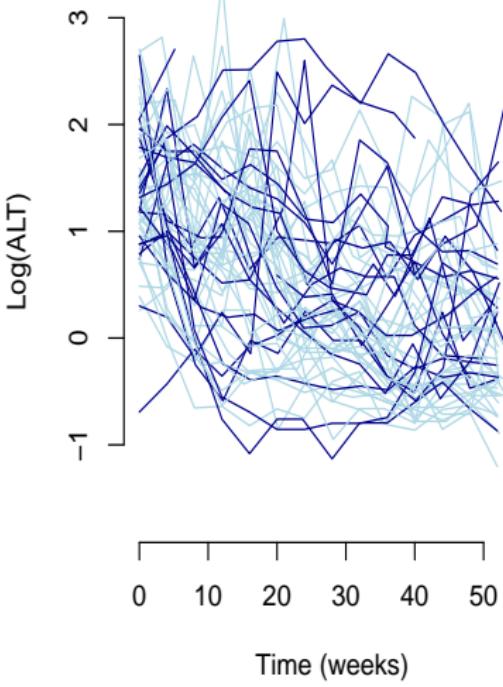
# Chronická hepatitida B

Shluky dle sdruženého modelu pro Log10(Viral load) a Log(ALT)

Model based group 0



Model based group 1



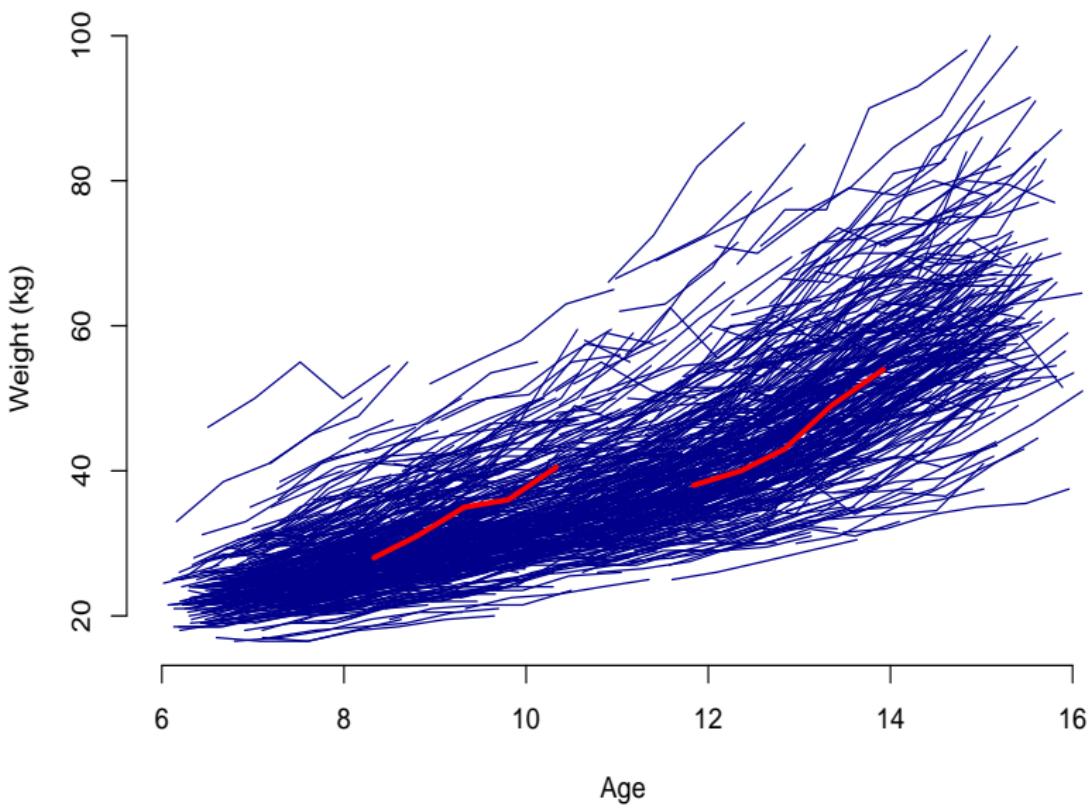
## Část IV

Stále v Králíkách

- ❖ Semilongitudinální studie (1997-2000)
- ❖ České děti ve věku 6–16 let
- ❖ Každé dítě měřeno/váženo až 6x v intervalu cca 6 měsíců
- ❖ Podrobnosti:  
Bláha, Krejčovský, Jiroutová, Kobzová, Sedlák, Brabec,  
Riedlová, Vignerová (2006).  
*Somatický vývoj současných českých dětí.*  
*Semilongitudinální studie.*  
UK v Praze a SZÚ Praha

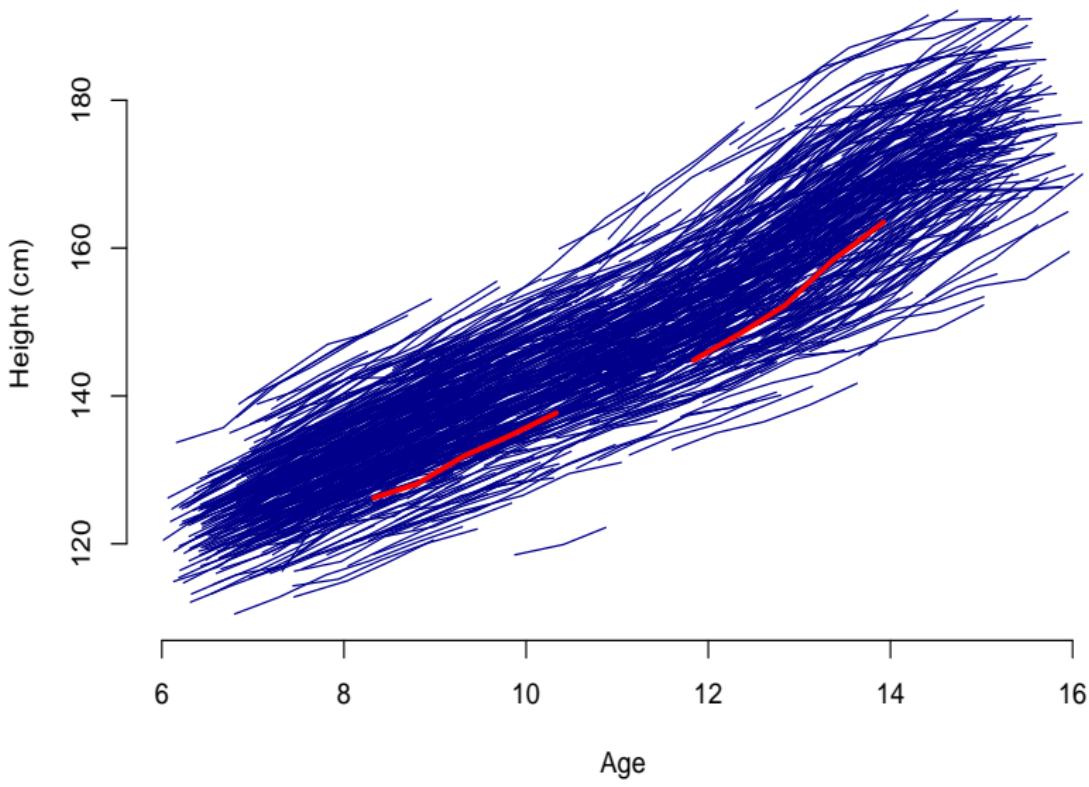
# Růst českých chlapců

Hmotnost



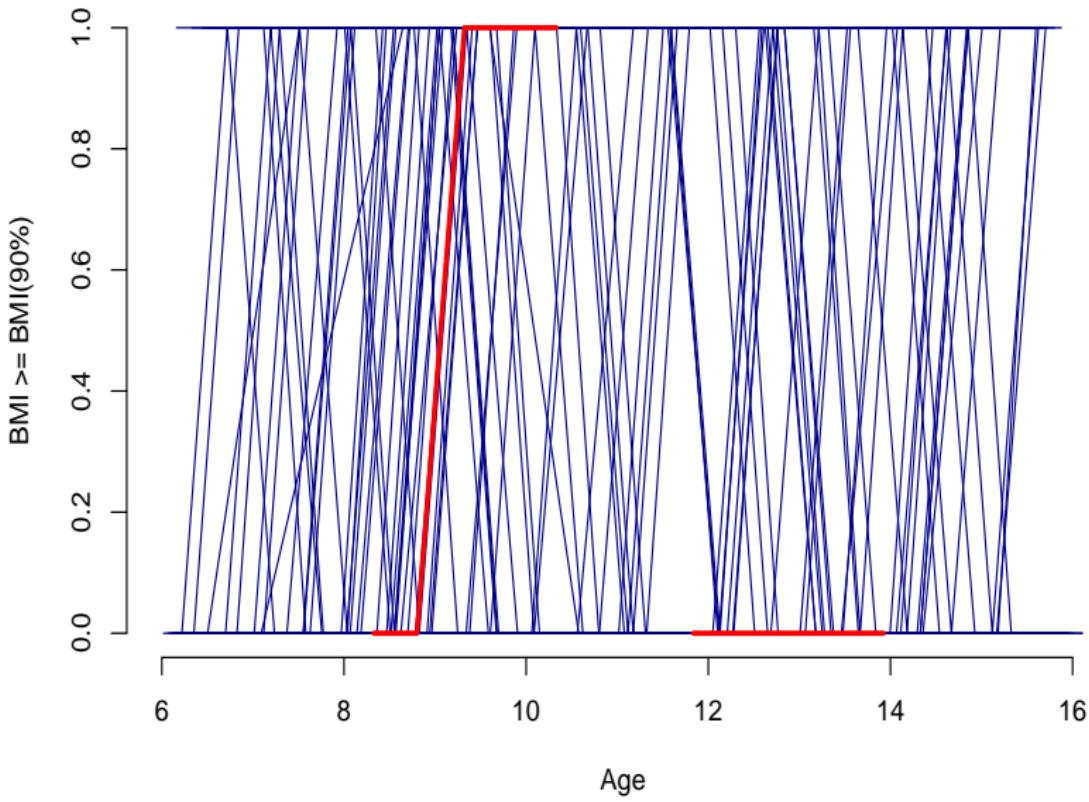
# Růst českých chlapců

Výška



# Růst českých chlapců

## BMI nad populačním 90% percentilem

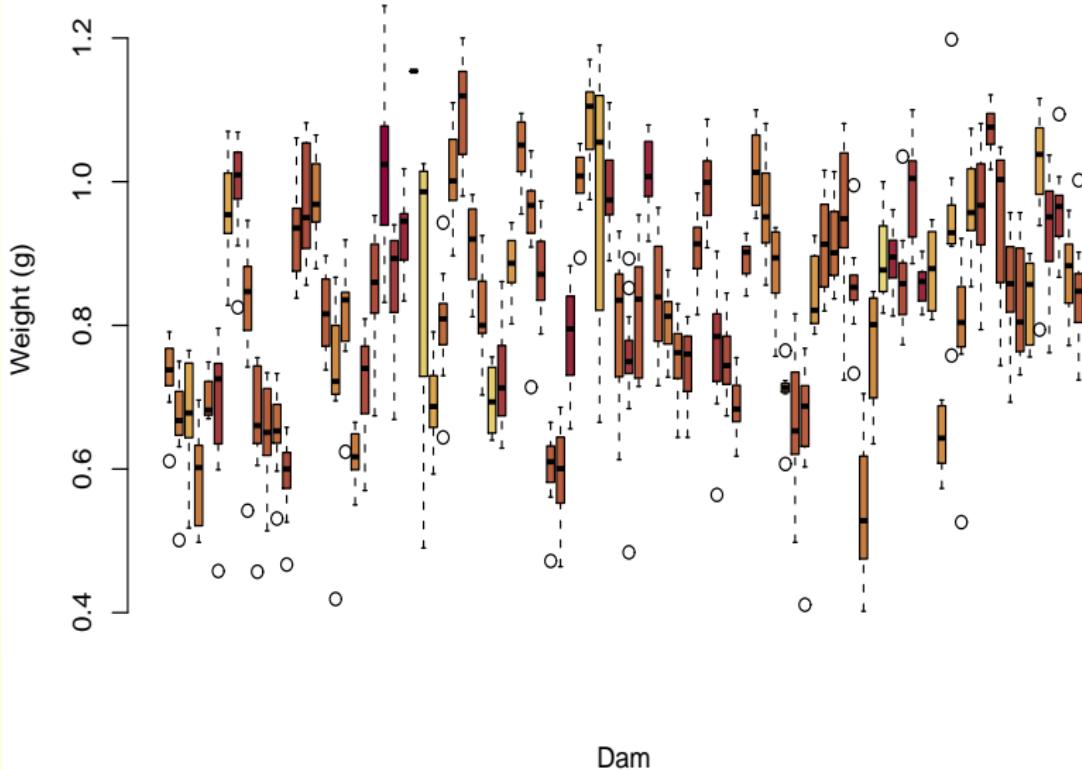


# Myši s etylen-glykolem

- ❖ Klinická studie (National Toxicology Program, National Institute of Health)
  - ❖ Těhotným myším byl podáván etylen-glykol
  - ❖ V 17. dnu gravidity byla zjištěna hmotnost plodů a fakt, je-li plod poškozen
  - ❖ Celkem 94 myší, každá 1–16 plodů
- 
- ❖ Ve statistické literatuře poměrně profláknutá data
    - ✳ Catalano, Ryan (1992, JASA)
    - ✳ :
    - ✳ Faes, Molenberghs, Aerts, Verbeke, Kenward (2009, Amer. Stat.)

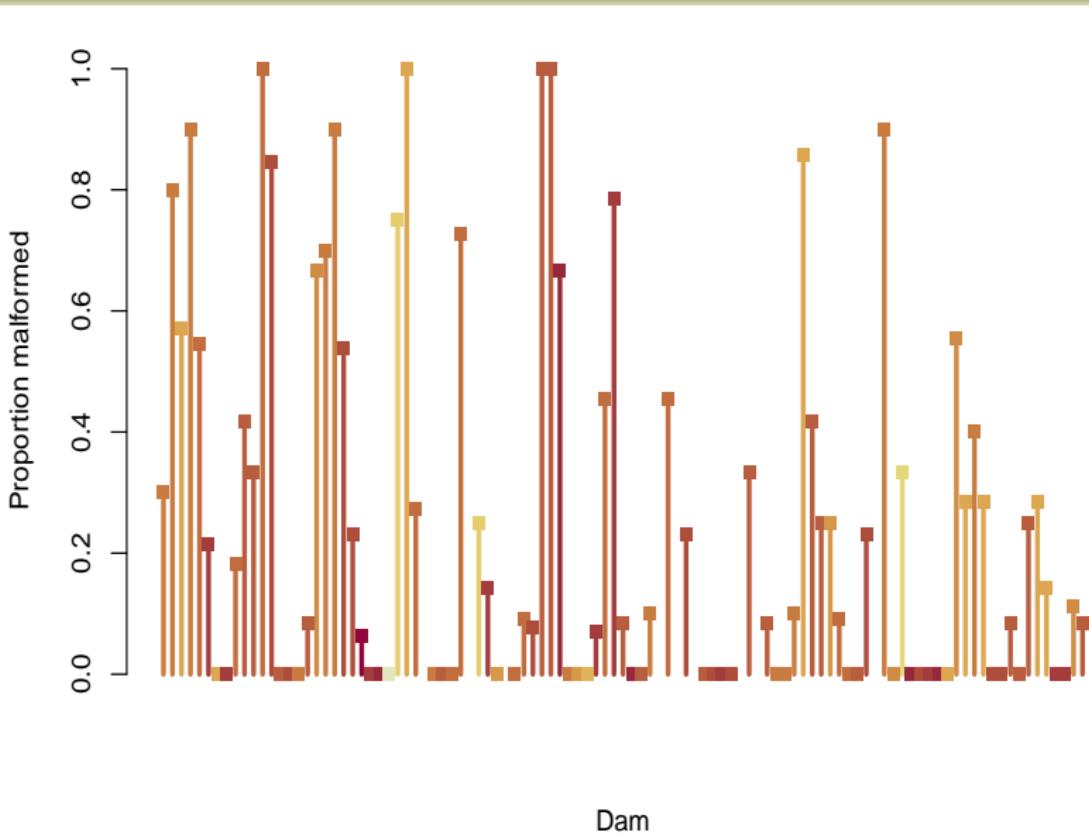
# Myši s etylen-glykolem

Hmotnost plodů



# Myši s etylen-glykolem

## Proporce poškozených plodů



# Myši s etylen-glykolem

- ❖ Nejsou to sice longitudinální data
- ❖ Ale jsou zde závislá pozorování (plody jedné samice)
- ❖ „čas“ ≡ index plodu v rámci samice
- Takže metodologicky se to od longitudinálních dat moc neliší

# Pravděpodobnostní model

- ❖ Jeden chlapec/jedna myš/jeden subjekt:

$$\mathbf{Y}_{i,1} = (Y_{i,1,1}, \dots, Y_{i,1,n_{i,1}})' \quad (\text{odezva č. 1})$$

$$\mathbf{Y}_{i,2} = (Y_{i,2,1}, \dots, Y_{i,2,n_{i,2}})' \quad (\text{odezva č. 2})$$

⋮

$$\mathbf{Y}_{i,R} = (Y_{i,R,1}, \dots, Y_{i,R,n_{i,R}})' \quad (\text{odezva č. R})$$

- ❖  $Y_{i,j,I}$  ne nutně spojité
- ❖ Budeme chtít shlukovat na základě

$$\mathbf{Y}_i = (\mathbf{Y}_{i,1}, \mathbf{Y}_{i,2}, \dots, \mathbf{Y}_{i,R})'$$

- ➡ Je potřeba mít **sdružený** model pro  $\mathbf{Y}_i$

# GLMM pro několik různých odezv

- ❖ Zobecněný lineární smíšený model (GLMM) pro  $r$ -tou odezvu ( $r = 1, \dots, R$ ):

$$h_r^{-1} \left\{ E(Y_{i,r,j} | \mathbf{b}_{i,r}) \right\} = \mathbf{x}'_{i,r,j} \boldsymbol{\alpha}_r + \mathbf{z}_{i,r,j} \mathbf{b}_{i,r}$$

- ✳  $h_r^{-1}$  ... linková funkce:
  - ✓ identita, je-li  $Y_{i,r} | \mathbf{b}_{i,r}$  normální
  - ✓ logit, je-li  $Y_{i,r} | \mathbf{b}_{i,r}$  alternativní (Bernoulli)
  - ✓ log, je-li  $Y_{i,r} | \mathbf{b}_{i,r}$  Poisson
- ✳  $\mathbf{x}_{i,1,1}, \dots, \mathbf{x}_{i,R,n_{i,R}}, \mathbf{z}_{i,1,1}, \dots, \mathbf{z}_{i,R,n_{i,R}}$   
... vektory regresorů
- ✳  $\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_R$  ... (neznámé) regresní parametry
- ✳  $\mathbf{b}_i = (\mathbf{b}'_{i,1}, \dots, \mathbf{b}'_{i,R})'$  i.i.d.  $\sum_{k=1}^K w_k \mathcal{N}(\boldsymbol{\mu}_k, \mathbb{D}_k)$

# Modelově založená shluková analýza

## Indikátory shluků

- ❖ Pro shlukování: i.i.d. náhodné veličiny  $U_1, \dots, U_n$
- ❖  $U_i \in \{1, \dots, K\}$
- ❖  $P(U_i = k) = w_k, k = 1, \dots, K, i = 1, \dots, n$

# Modelově založená shluková analýza

Vícerozměrný GLMM s normální směsí v rozdělení náhodných efektů zapsaný hierarchicky

- ❖ Hierarchicky zapsáno:

$$\underbrace{[\mathbf{Y}_{i,1}, \dots, \mathbf{Y}_{i,R}] \mid \mathbf{b}_{i,1}, \dots, \mathbf{b}_{i,R}}_{\mathbf{Y}_i} \sim \prod_{r=1}^R \prod_{j=1}^{n_{i,r}} \underbrace{[y_{i,r,j} \mid \mathbf{b}_{i,r}]}_{\text{normální, Bernoulli, Poisson, ... se střední hodnotou } h_r(\mathbf{x}_{i,r,j}\boldsymbol{\alpha}_r + \mathbf{z}_{i,r,j}\mathbf{b}_{i,r})}$$

$$\left. \begin{array}{lcl} P(U_i = k) & = & w_k \\ [\mathbf{b}_i \mid U_i = k] & \sim & \mathcal{N}(\boldsymbol{\mu}_k, \mathbb{D}_k) \end{array} \right\} \quad k = 1, \dots, K$$

# Modelově založená shluková analýza

Vícerozměrný GLMM s normální směsí v rozdělení náhodných efektů

$$\Rightarrow p(\mathbf{y}_i) = \int_{\mathbb{R}^q} \left\{ \prod_{r=1}^R \prod_{j=1}^{n_{i,r}} \underbrace{p(y_{i,r,j} | \mathbf{b}_{i,r})}_{\text{normální, Bernoulli, Poisson, ...}} \right\} \times$$

$$\underbrace{\left\{ \sum_{k=1}^K w_k \varphi(\mathbf{b}_i; \mu_k, \mathbb{D}_k) \right\}}_{\text{normální směs}} d\mathbf{b}_i$$

$$\Rightarrow p(\mathbf{y}_i | U_i = k) =$$

$$\int_{\mathbb{R}^q} \left\{ \prod_{r=1}^R \prod_{j=1}^{n_{i,r}} \underbrace{p(y_{i,r,j} | \mathbf{b}_{i,r})}_{\text{normální, Bernoulli, Poisson, ...}} \right\} \times \underbrace{\varphi(\mathbf{b}_i; \mu_k, \mathbb{D}_k)}_{\text{normální}} d\mathbf{b}_i$$

# Modelově založená shluková analýza

Vícerozměrný GLMM s normální směsí v rozdělení náhodných efektů

## ❖ Model:

$$[\mathbf{Y}_i | \mathbf{b}_i] \sim \prod_{r=1}^R \prod_{j=1}^{n_{i,r}} \text{GLMM}(Y_{i,r,j} | \mathbf{x}'_{i,r,j} \boldsymbol{\alpha}_r + \mathbf{z}'_{i,r,j} \mathbf{b}_{i,r})$$

$$\left. \begin{array}{lcl} \mathsf{P}(U_i = k) & = & w_k \\ [\mathbf{b}_i | U_i = k] & \sim & \mathcal{N}(\mu_k, \mathbb{D}_k) \end{array} \right\} \quad k = 1, \dots, K$$

## ❖ Parametry modelu:

$$\theta = (\boldsymbol{\alpha}', \mathbf{w}', \mu'_1, \dots, \mu'_K, \text{vec}(\mathbb{D}_1)', \dots, \text{vec}(\mathbb{D}_K)', \text{rozptyl GLMM})' \in \Theta$$

## ❖ Nepřímo pozorovatelné náhodné veličiny/vektory:

$$\eta = (\mathbf{b}'_1, \dots, \mathbf{b}'_n, U_1, \dots, U_n)' \in \mathbb{R}^q \times \dots \times \mathbb{R}^q \times \{1, \dots, K\}^n$$

# Modelově založená shluková analýza

Vícerozměrný GLMM s normální směsí v rozdělení náhodných efektů

➡ Bayesova věta:

$$\pi_{i,k}(\theta) \equiv P(U_i = k \mid \mathbf{y}_i, \theta)$$

$$\propto w_k p(\mathbf{y}_i \mid U_i = k, \theta)$$

$$= w_k \int_{\mathbb{R}^q} \left\{ \prod_{r=1}^R \prod_{j=1}^{n_{i,r}} \underbrace{p(y_{i,r,j} \mid \mathbf{b}_{i,r})}_{\text{normální, Bernoulli, Poisson, ...}} \right\} \times$$

$$\underbrace{\varphi(\mathbf{b}_i; \mu_k, \mathbb{D}_k)}_{\text{normální}} d\mathbf{b}_i$$

- ☞ I při známé hodnotě  $\theta$  analyticky obecně nevyjádřitelné
- ☞ Navíc,  $\theta$  je potřeba nejprve odhadnout...

## Část V

Co v Králíkách bude možná až po 22. hodině

# Vícerozměrný GLMM s normální směsí

## Odhad parametrů

- ❖ Parametry modelu:

$$\boldsymbol{\theta} = (\boldsymbol{\alpha}', \mathbf{w}', \mu_1', \dots, \mu_K', \\ \text{vec}(\mathbb{D}_1)', \dots, \text{vec}(\mathbb{D}_K)', \text{rozptyl GLMM})' \in \Theta$$

- ❖ Věrohodnost:

$$L(\boldsymbol{\theta}) = \prod_{i=1}^n p(\mathbf{y}_i | \boldsymbol{\theta}) = \prod_{i=1}^n \int_{\mathbb{R}^q} \left\{ \prod_{r=1}^R \prod_{j=1}^{n_{i,r}} \underbrace{p(y_{i,r,j} | \boldsymbol{\alpha}, \mathbf{b}_{i,r})}_{\text{normální, Bernoulli, Poisson, ...}} \right\} \times$$

$$\underbrace{\left\{ \sum_{k=1}^K w_k \varphi(\mathbf{b}_i; \mu_k, \mathbb{D}_k) \right\} d\mathbf{b}_i}_{\text{normální směs}}$$

# Vícerozměrný GLMM s normální směsí

Odhad parametrů, maximální věrohodnost

- ❖  $\max_{\theta \in \Theta} \log\{L(\theta)\}$  by šlo hledat EM algoritmem
- ❖  $\theta^{(m)}$  ... hodnota  $\theta$  v  $m$ -té iteraci EM algoritmu

E-krok: spočítej (pro  $i = 1, \dots, n, k = 1, \dots, K$ )

$$\pi_{i,k}(\theta^{(m)}) \equiv P(U_i = k \mid \mathbf{y}_i, \theta^{(m)})$$

$$\propto w_k^{(m)} \int_{\mathbb{R}^q} \left\{ \prod_{r=1}^R \underbrace{\prod_{j=1}^{n_{i,r}} p(y_{i,r,j} \mid \alpha^{(m)}, \mathbf{b}_{i,r})}_{\text{normální, Bernoulli, Poisson, ...}} \right\} \times \underbrace{\varphi(\mathbf{b}_i; \mu_k^{(m)}, \mathbb{D}_k^{(m)})}_{\text{normální}} d\mathbf{b}_i$$

# Vícerozměrný GLMM s normální směsí

Odhad parametrů, maximální věrohodnost

- ❖  $\max_{\theta \in \Theta} \log\{L(\theta)\}$  by šlo hledat EM algoritmem
- ❖  $\theta^{(m)}$  ... hodnota  $\theta$  v  $m$ -té iteraci EM algoritmu

M-krok: maximalizuj vzhledem k  $\theta$

$$Q(\theta | \theta^{(m)}) = \sum_{i=1}^n \sum_{k=1}^K \pi_{i,k}(\theta^{(m)}) \log w_k +$$
$$\sum_{i=1}^n \sum_{k=1}^K \pi_{i,k}(\theta^{(m)}) \log \left[ \int_{\mathbb{R}^q} \left\{ \prod_{r=1}^R \prod_{j=1}^{n_{i,r}} \underbrace{p(y_{i,r,j} | \alpha, \mathbf{b}_{i,r})}_{\text{normální, Bernoulli, Poisson, ...}} \right\} \times \right.$$
$$\left. \underbrace{\varphi(\mathbf{b}_i; \mu_k, \mathbb{D}_k) d\mathbf{b}_i}_{\text{normální}} \right]$$

# Vícerozměrný GLMM s normální směsí

Bayesovský odhad parametrů

## ❖ Parametry modelu:

$$\boldsymbol{\theta} = (\boldsymbol{\alpha}', \mathbf{w}', \mu_1', \dots, \mu_K',$$

$$\text{vec}(\mathbb{D}_1)', \dots, \text{vec}(\mathbb{D}_K)', \text{rozptyl GLMM})' \in \Theta$$

## ❖ Rušivé parametry:

$$\boldsymbol{\eta} = (\mathbf{b}_1', \dots, \mathbf{b}_n', \mathbf{U}_1, \dots, \mathbf{U}_n)' \in \mathbb{R}^q \times \dots \times \mathbb{R}^q \times \{1, \dots, K\}^n$$

## ❖ Věrohodnost:

$$L(\boldsymbol{\theta}, \boldsymbol{\eta}) = \prod_{i=1}^n p(\mathbf{y}_i | \boldsymbol{\theta}, \boldsymbol{\eta}) = \prod_{i=1}^n \prod_{r=1}^R \prod_{j=1}^{n_{i,r}} \underbrace{p(y_{i,r,j} | \boldsymbol{\alpha}, \mathbf{b}_{i,r})}_{\text{normální, Bernoulli, Poisson, ...}}$$

☞ Věrohodnost „obyčejného“ GLM (pouze jednou **M**)

# Vícerozměrný GLMM s normální směsí

Bayesovský odhad parametrů

❖ Apriorní rozdělení:

$$p(\theta, \eta) = p(\eta | \theta) \times p(\theta)$$

$$= \prod_{i=1}^n \left\{ p(\mathbf{b}_i | U_i, \theta) \times p(U_i | \theta) \right\} \times p(\theta)$$

- \*  $p(\mathbf{b}_i | U_i = k, \theta) \sim \mathcal{N}(\mu_k, \mathbb{D}_k)$
- \*  $p(U_i | \theta) \sim P(U_i = k | \mathbf{w}) = w_k$
- \*  $p(\theta) \equiv$  „slušné“ apriorní rozdělení  
pro „klasické“ parametry

# Vícerozměrný GLMM s normální směsí

Bayesovský odhad parametrů

- ❖ Aposteriorní rozdělení:

$$p(\theta, \eta | \mathbf{y}) = \frac{L(\theta, \eta) \times p(\theta, \eta)}{\int L(\theta, \eta) \times p(\theta, \eta) d\theta d\eta}$$

- ❖ MCMC (detaily pro zájemce dnes po 22. hodině):

⇒ Náhodný výběr  $(\theta^{(1)}, \eta^{(1)}), \dots, (\theta^{(M)}, \eta^{(M)})$  z  $p(\theta, \eta | \mathbf{y})$

# Modelově založená shluková analýza

Shlukování na základě MCMC

- ❖ Připomenutí: shlukujeme na základě

$$P(U_i = k | \mathbf{y}) = E_{\theta}\{P(U_i = k | \mathbf{y}, \theta) | \mathbf{y}\} = E_{\theta}\{\pi_{i,k}(\theta) | \mathbf{y}\}$$

- ❖ MCMC odhad:

$$\widehat{P}(U_i = k | \mathbf{y}) = M^{-1} \sum_{m=1}^M \pi_{i,k}(\theta^{(m)})$$

- ❖ Bylo:

$$\pi_{i,k}(\theta^{(m)}) \propto$$

$$w_k^{(m)} \int_{\mathbb{R}^q} \left\{ \prod_{r=1}^R \prod_{j=1}^{n_{i,r}} \underbrace{p(y_{i,r,j} | \alpha^{(m)}, \mathbf{b}_{i,r})}_{\text{normální, Bernoulli, Poisson, ...}} \right\} \times \underbrace{\varphi(\mathbf{b}_i; \mu_k^{(m)}, \mathbb{D}_k^{(m)})}_{\text{normální}} d\mathbf{b}_i$$

# Modelově založená shluková analýza

Shlukování na základě MCMC, alternativně

- ❖ Připomenutí: shlukujeme na základě

$$\text{P}(U_i = k \mid \mathbf{y}) = \text{E}_{\theta, \mathbf{b}_i} \{ \text{P}(U_i = k \mid \mathbf{y}, \theta, \mathbf{b}_i) \mid \mathbf{y} \}$$

- ❖ MCMC odhad:

$$\widehat{\text{P}}(U_i = k \mid \mathbf{y}) = M^{-1} \sum_{m=1}^M \text{P}(U_i = k \mid \mathbf{y}, \theta^{(m)}, \mathbf{b}_i^{(m)})$$

- ❖ Trocha počítání:

$$\text{P}(U_i = k \mid \mathbf{y}, \theta^{(m)}, \mathbf{b}_i^{(m)}) \propto w_k^{(m)} \prod_{r=1}^R \underbrace{\prod_{j=1}^{n_{i,r}} p(y_{i,r,j} \mid \alpha^{(m)}, \mathbf{b}_{i,r})}_{\text{normální, Bernoulli, Poisson, ...}}$$

# Modelově založená shluková analýza

Shlukování na základě MCMC, další alternativa

- ❖ Připomenutí: shlukujeme na základě

$$\mathsf{P}(U_i = k \mid \mathbf{y}) = \mathsf{P}(V_{i,k} = 1 \mid \mathbf{y})$$

- \*  $\mathbf{V}_i = (V_{i,1}, \dots, V_{i,K})'$ ,  $V_{i,k} \in \{0, 1\}$ ,  $\sum_{k=1}^K V_{i,k} = 1$
- \* Jinými slovy:  $V_{i,k} = \mathbb{I}(U_i = k)$

- ❖ MCMC odhad:

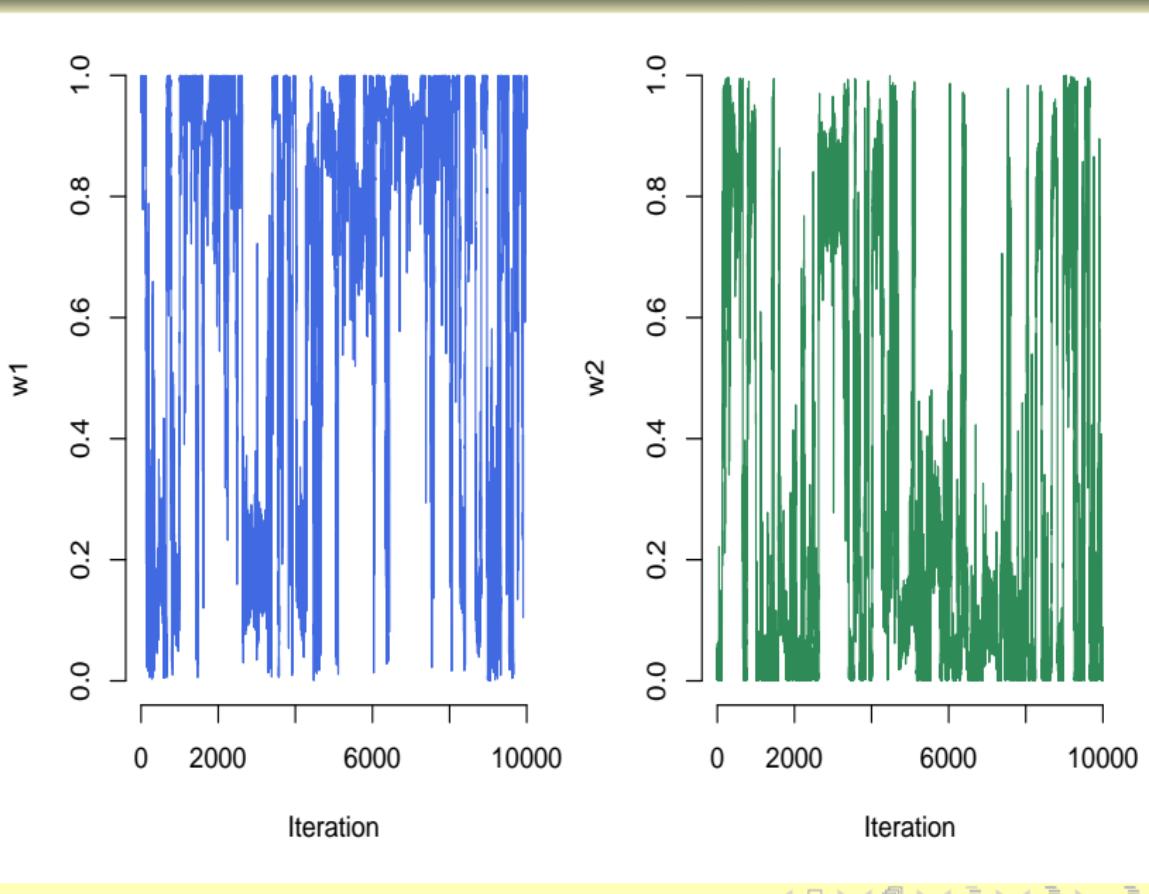
$$\widehat{\mathsf{P}}(U_i = k \mid \mathbf{y}) = M^{-1} \sum_{m=1}^M \mathbb{I}(U_i^{(m)} = k)$$

# Modelově založená shluková analýza

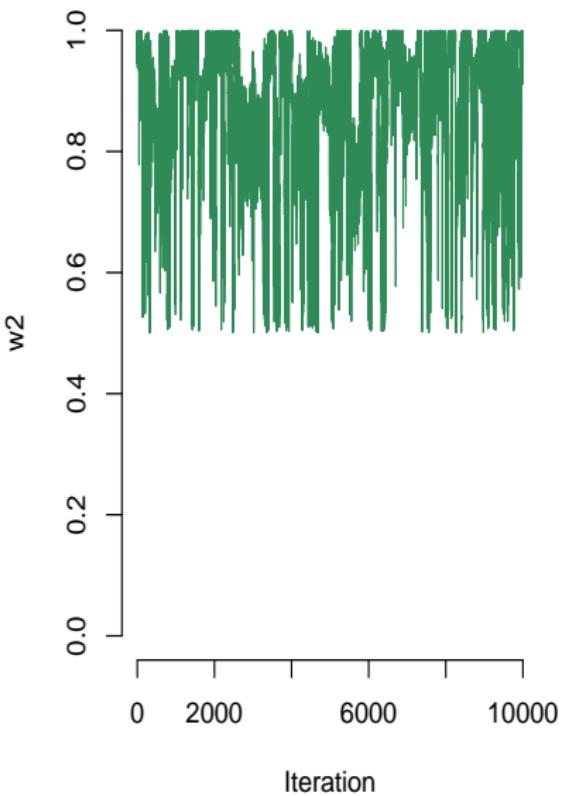
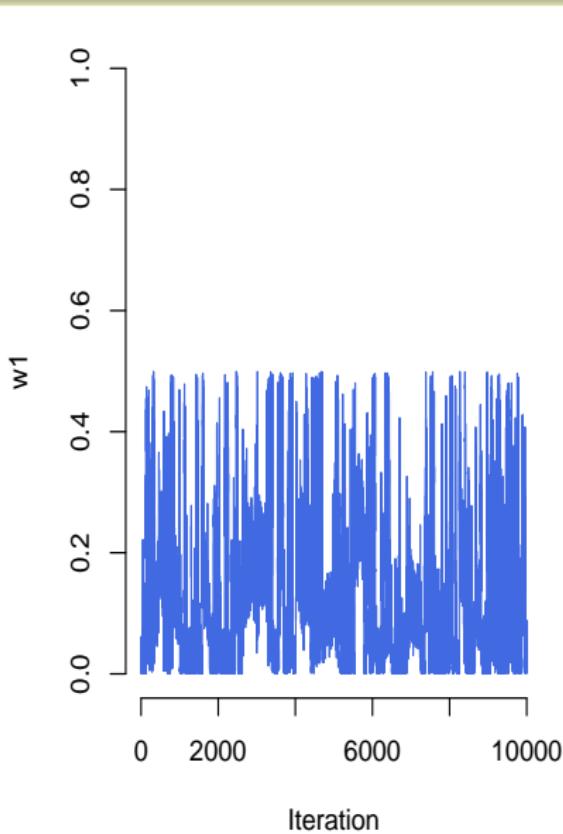
## Problémy se shlukováním na základě MCMC

- ❖ Aposteriorní rozdělení je invariantní vůči změně očíslování komponent
- ❖ Existuje  $K!$  různých očíslování
- ➡ Aposteriorní rozdělení má  $K!$  symetrických lépe či hůře separovaných nosičů
- ☞ Je-li  $(\theta^{(m)}, \eta^{(m)})$  stav MCMC v  $m$ -té iteraci, je potřeba vědět, ke kterému z  $K!$  nosičů patří
- ☞ Jednoduchá zjednoznačnění aposteriorního rozdělení pomocí omezení typu  $w_1 < \dots < w_K$  nelze v rámci MCMC použít, jsou-li jednotlivé nosiče blízko sebe

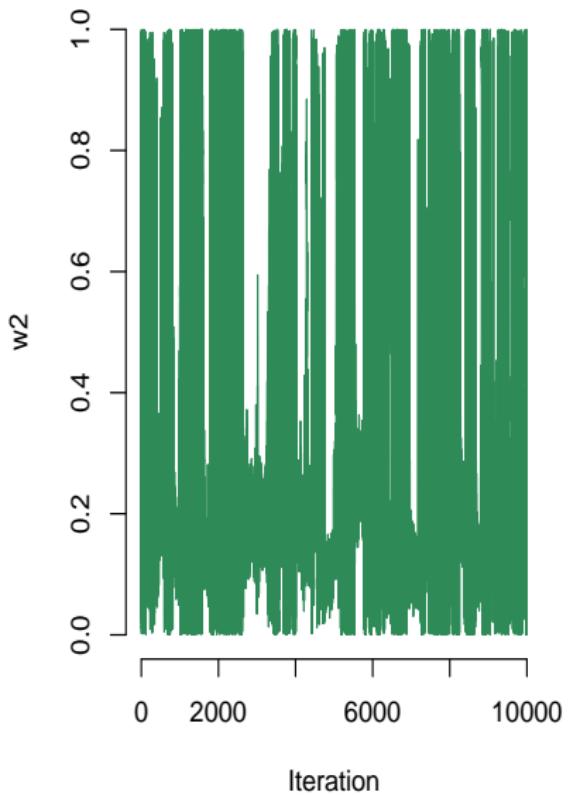
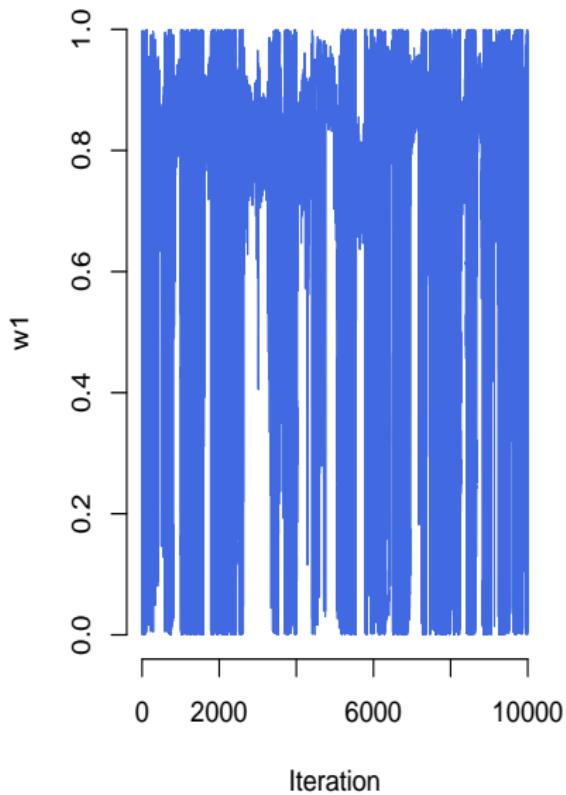
# K = 2, MCMC bez omezení



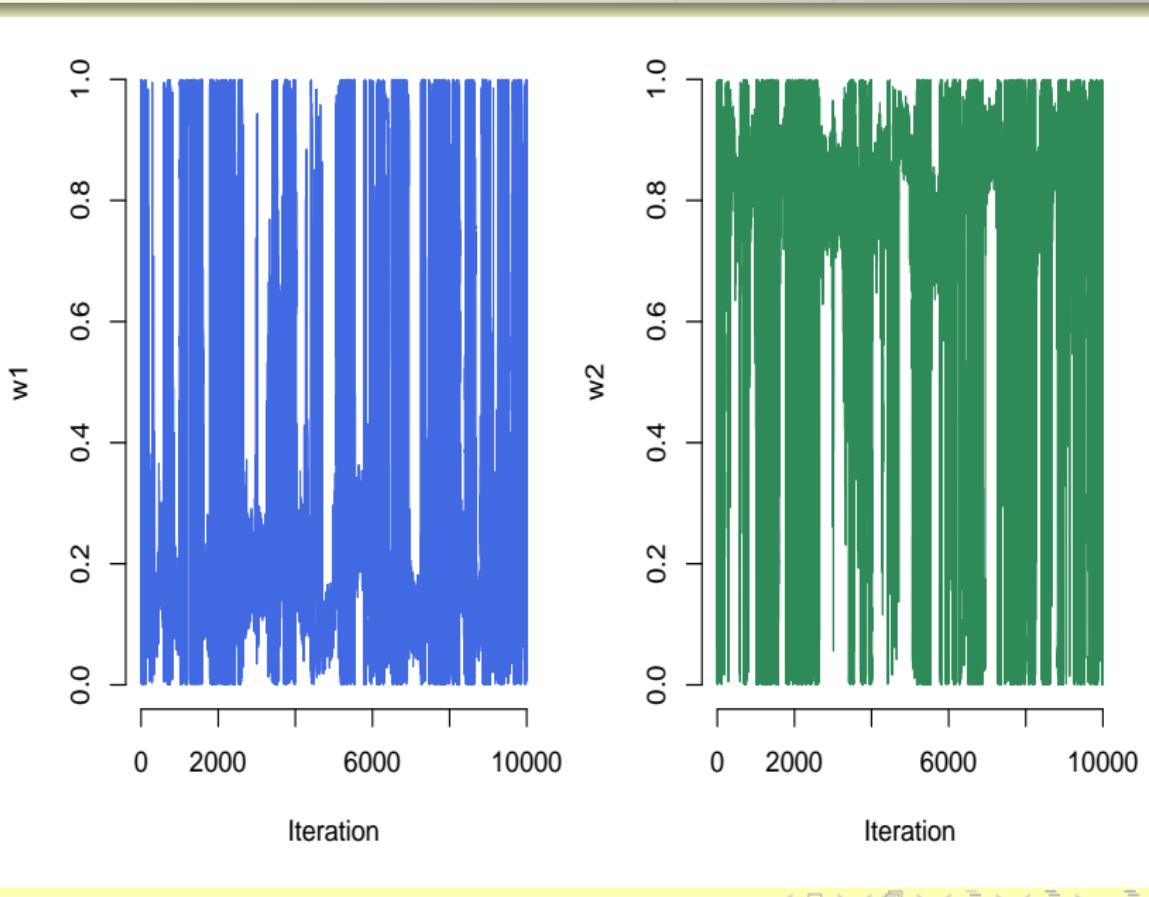
# $K = 2$ , identifikační omezení $w_1 < w_2$



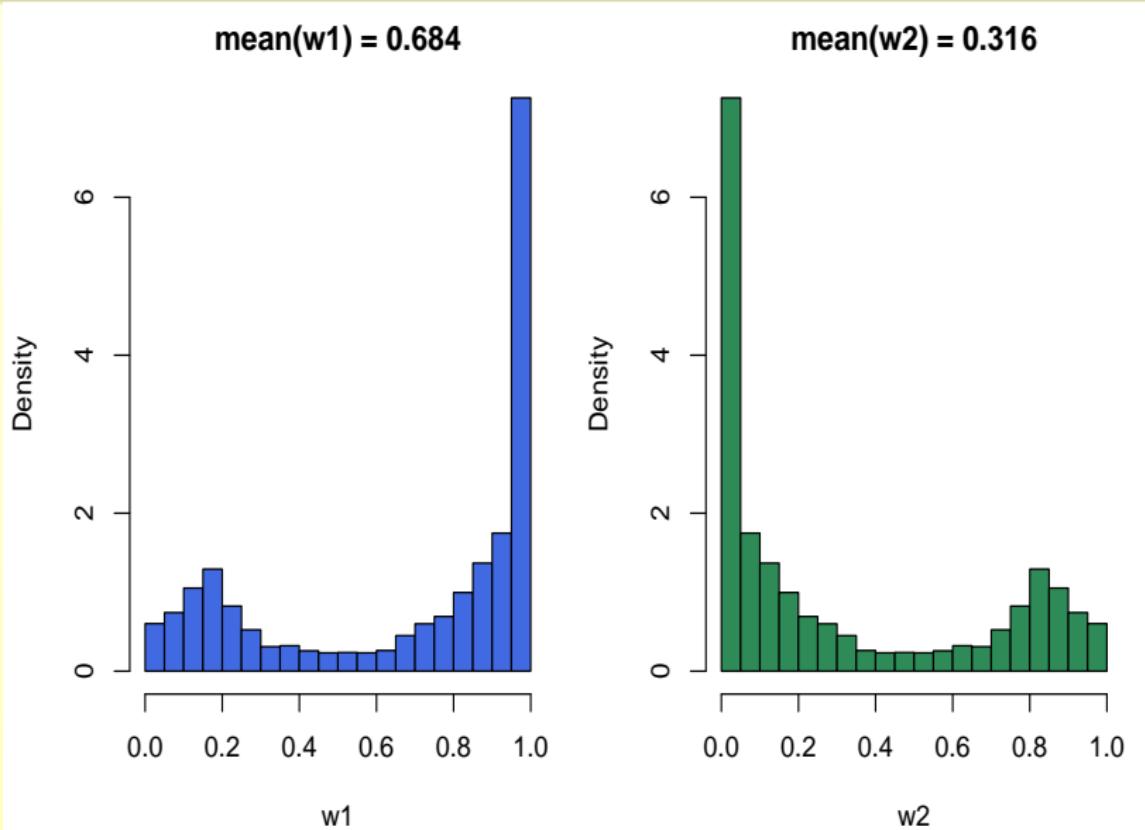
# $K = 2$ , identifikační omezení $\mu_{1,1} < \mu_{2,1}$



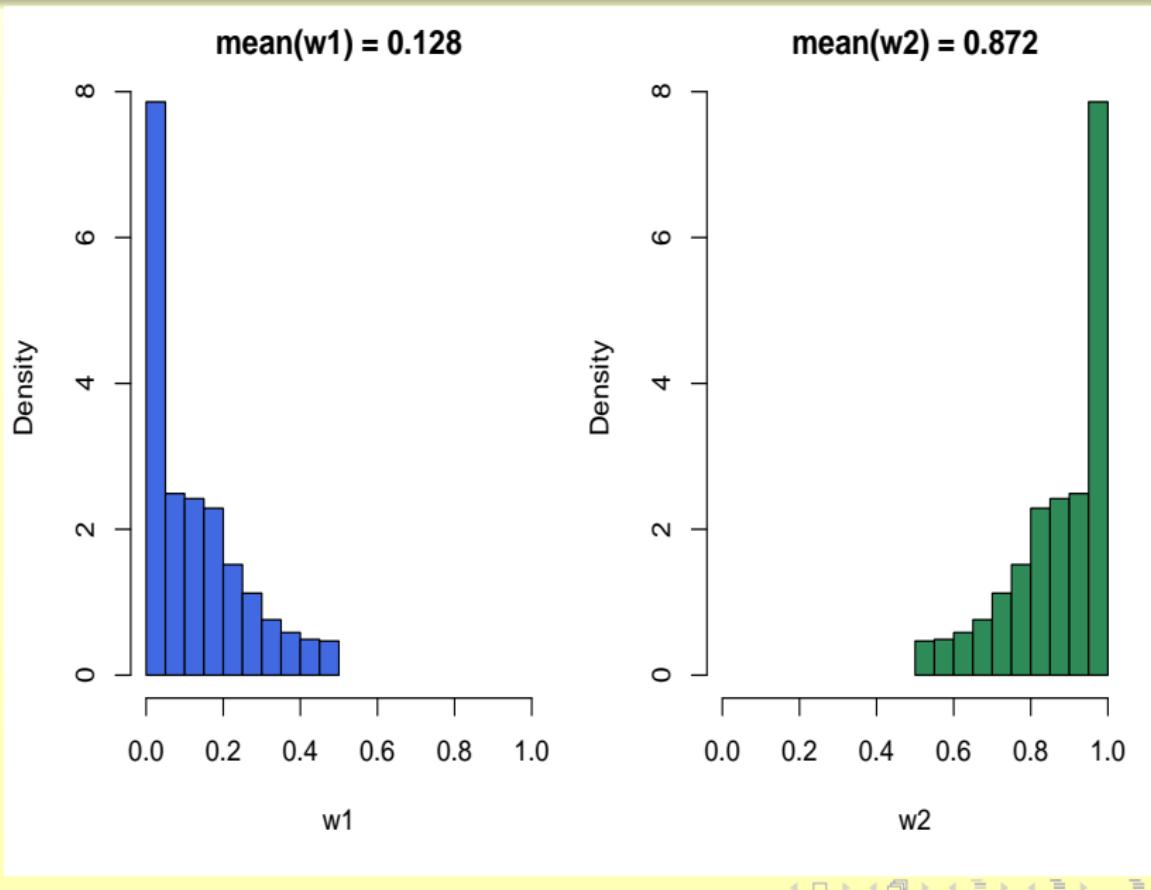
# $K = 2$ , identifikační omezení $\mu_{1,2} < \mu_{2,2}$



# K = 2, MCMC bez omezení

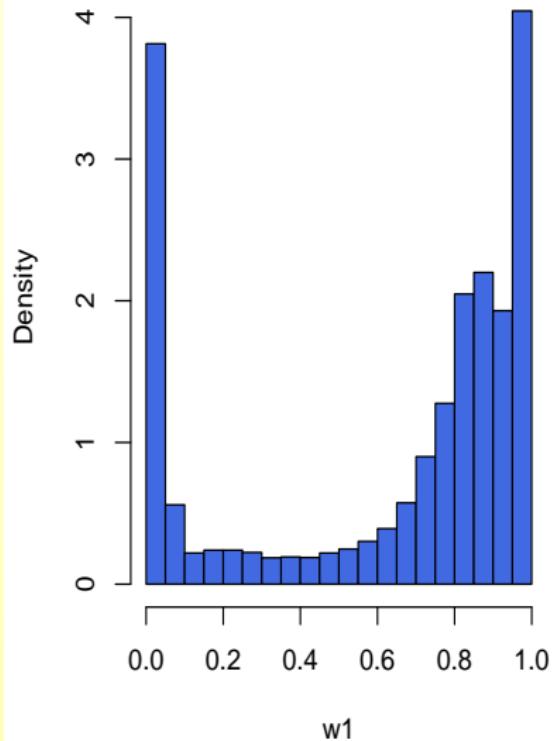


# $K = 2$ , identifikační omezení $w_1 < w_2$

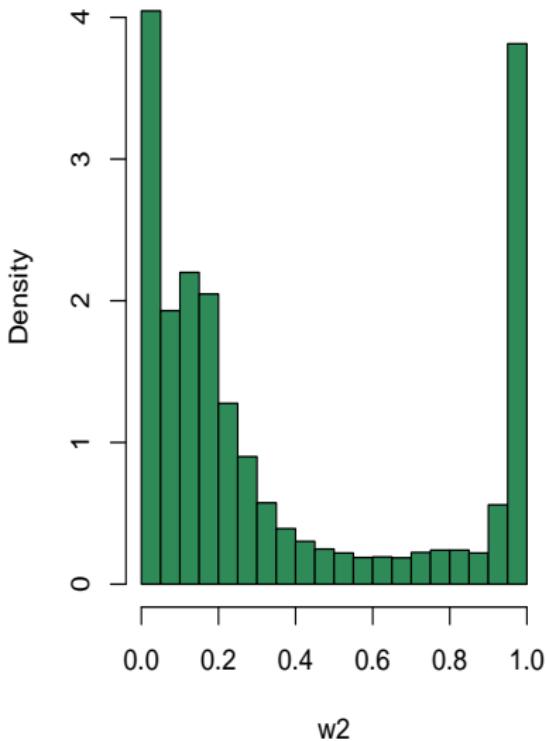


$K = 2$ , identifikační omezení  $\mu_{1,1} < \mu_{2,1}$

$\text{mean}(w1) = 0.628$



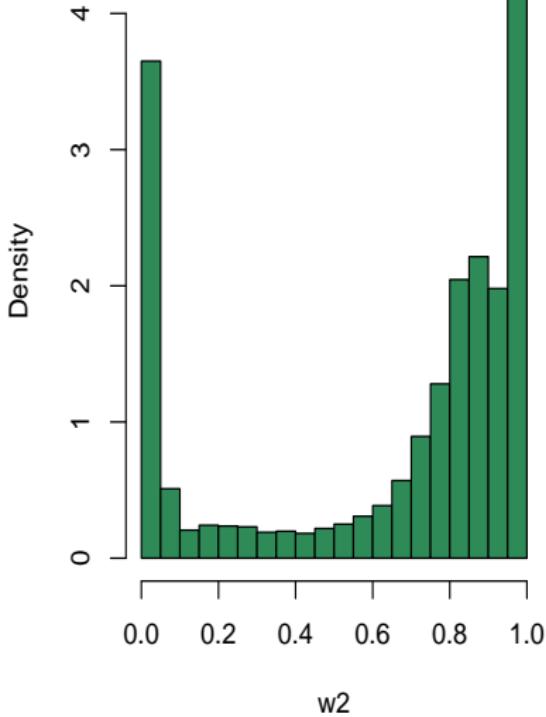
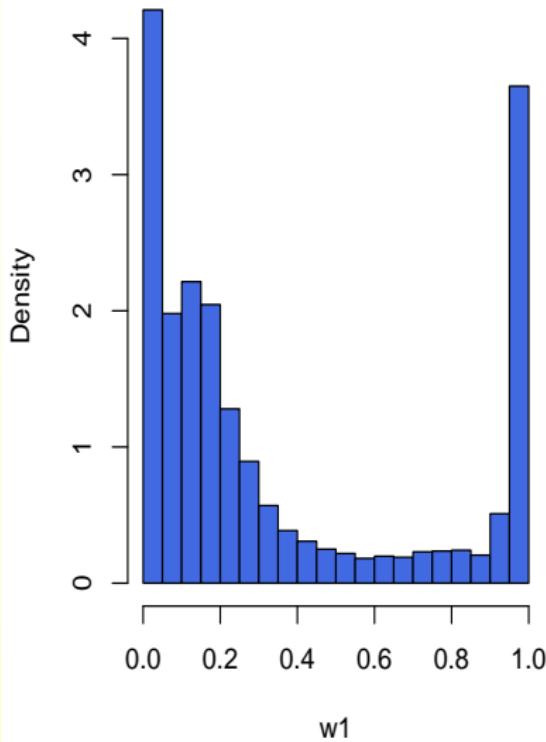
$\text{mean}(w2) = 0.372$



$K = 2$ , identifikační omezení  $\mu_{1,2} < \mu_{2,2}$

$\text{mean}(w1) = 0.362$

$\text{mean}(w2) = 0.638$



# Modelově založená shluková analýza

## Problémy se shlukováním na základě MCMC

- ❖ Aby šel výstup z MCMC použít ke shlukování, je potřeba použít sofistikovanější přístup k identifikaci jednoho z  $K!$  nosičů
- ❖ Souhrnně na toto téma:  
[Jasra, Holmes, Stephens \(2005\).](#)  
Markov chain Monte Carlo methods and the label switching problem in Bayesian mixture modeling.  
*Statistical Science*, **20**, 50–67
- ❖ Se mnou na toto téma + aplikace na české děti  
či myši ošetřené etylen-glykolem:  
Zase někdy jindy...
  - ☞ Třeba na PASTi ve středu odpoledne v Karlíně
  - ☞ Nebo v Řecku v srpnu na EMS

## Část VI

Něco na závěr

- ❖ R balíček mixAK
  - ✿ Byla o něm řeč na XV. ROBUSTu v Roháčích
  - ✿ Tenkrát tohle všechno ještě neuměl
- ❖ <http://cran.R-project.org/package=mixAK>

## Za data děkuji:

**Právníci:** doc. K. Zvára (MFF UK), prof. D. Hendrych (PF UK)

**Hepatitida B:** H. Janssen, B. Hansen (Erasmus Medisch Centrum Rotterdam)

**České děti:** ing. J. Vignerová, dr. M. Brabec (SZÚ)

**Myši s etylen-glykolem:** National Toxicology Program (<http://ntp.niehs.nih.gov>), C. Faes (Univ. Hasselt)

## Za finance děkuji:

GAČRu 201/09/P077

MŠMT ČR MSM 0021620839

Děkuji za pozornost!

## Reklama

25th International Workshop on Statistical Modelling

4.–9. července/júla 2010

Glasgow

<http://mathstore.ac.uk/iwsm/>

- Jednodenní krátký kurz:  
Giles Hooker – Analýza funkcionálních dat
- Studentská soutěž
- Abstrakty do 5. února/februára 2010 (zítra ...)
- ☞ Alternativní program místo dnešního večírku???