

Katedra pravděpodobnosti a matematické statistiky



MATEMATICKO-FYZIKÁLNÍ  
FAKULTA

Univerzita Karlova

---

doc. RNDr. Arnošt Komárek, Ph.D.

**NMST431 Bayesovské metody**

---

Zimní semestr 2021–22

# 1

## Lineární model

$$\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \mathbb{X} : \text{pevná matice } n \times k,$$
$$\boldsymbol{\varepsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n).$$

⇒  $\mathbf{Y} \sim \mathcal{N}_n(\mathbb{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$

- **Parametry:**  $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \tau)^\top$ , kde  $\tau = \sigma^{-2} > 0$

## Příklad: Vážení lehkých objektů

- Potřeba zjistit hmotnost dvou velice lehkých objektů.
  - $\beta_1, \beta_2$ : hmotnost objektu A, resp. B
- Pokus (naměřené hodnoty v  $\mu\text{g}$ )
  - ▮  $Y_i, i = 1, \dots, n$  ( $n = 18$ ):
    - 2× zvážen objekt A: 109, 85
    - 9× zvážen objekt B: 114, 121, 140, 122, 125, 129, 98, 134, 133
    - 7× zváženy oba objekty A a B: 217, 203, 243, 229, 233, 221, 221
- Každé měření má (náhodnou) chybu, o které budeme předpokládat, že má normální rozdělení s nulovou střední hodnotou a rozptylem  $\sigma^2$ .
- Lineární model pro  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ :
$$\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\beta} = (\beta_1, \beta_2)^\top$$
  - Jak vypadá matice  $\mathbb{X}$ ?

$$\begin{aligned}L(\boldsymbol{\theta}) &= p(\mathbf{y} | \boldsymbol{\theta}) = \dots \\&= (2\pi)^{-\frac{n}{2}} |\mathbb{X}^T \mathbb{X}|^{\frac{1}{2}} \tau^{\frac{n}{2}} \exp \left[ -\frac{\tau}{2} \left\{ \text{SS}_e + (\boldsymbol{\beta} - \mathbf{b})^T \mathbb{X}^T \mathbb{X} (\boldsymbol{\beta} - \mathbf{b}) \right\} \right] \\&= (2\pi)^{-\frac{n}{2}} |\mathbb{X}^T \mathbb{X}|^{\frac{1}{2}} \tau^{\frac{k}{2}} \exp \left\{ -\frac{\tau}{2} (\boldsymbol{\beta} - \mathbf{b})^T \mathbb{X}^T \mathbb{X} (\boldsymbol{\beta} - \mathbf{b}) \right\} \tau^{\frac{n-k+2}{2}-1} \exp \left( -\tau \frac{\text{SS}_e}{2} \right)\end{aligned}$$

- kde
- $\mathbf{b} = \mathbf{b}(\mathbf{y}) = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbf{y}$
  - $\text{SS}_e = \text{SS}_e(\mathbf{y}) = (\mathbf{y} - \mathbb{X}\mathbf{b})^T (\mathbf{y} - \mathbb{X}\mathbf{b})$
  - $s = s(\mathbf{y}) = \sqrt{\frac{\text{SS}_e}{n-k}}$

**Příklad: Lehké objekty** •  $\mathbf{b} = (98.89, 124.42)$

•  $\text{SS}_e = 2525.7$

•  $s = 12.56$

# Neinformativní apriorní rozdělení

---

- $p(\beta) \propto 1, \beta \in \mathbb{R}^k$

- $p(\sigma^2) \propto \frac{1}{\sigma^2}, \sigma^2 > 0$

- odpovídá  $p(\log \sigma^2) \propto c$ , resp.  $p(\log \sigma) \propto c$

- $p(\tau) \propto \frac{1}{\tau}, \tau > 0$

- $\beta$  a  $\tau$  apriori nezávislé

$$p(\boldsymbol{\beta}, \tau | \mathbf{y}) \propto \dots \propto$$

$$\tau^{\frac{n}{2}-1} \exp\left[-\frac{\tau}{2}\left\{\mathbf{S}\mathbf{S}_e + (\boldsymbol{\beta} - \mathbf{b})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \mathbf{b})\right\}\right], \quad \tau > 0$$

- Za jakých podmínek se jedná skutečně o (nedegenerované) rozdělení?
- To jest, za jakých podmínek existuje integrál z výše uvedené funkce?

$$p(\boldsymbol{\beta}, \tau | \mathbf{y}) = p(\boldsymbol{\beta} | \tau, \mathbf{y}) \times p(\tau | \mathbf{y})$$

$$p(\beta | \tau, \mathbf{y}) \propto \dots \propto \exp \left[ -\frac{\tau}{2} \left\{ (\beta - \mathbf{b})^\top \mathbb{X}^\top \mathbb{X} (\beta - \mathbf{b}) \right\} \right]$$

- ▶ Normující konstanta je proporcionální determinantu matice  $|\mathbb{X}^\top \mathbb{X}|^{-\frac{1}{2}}$
- ▶ To jest, aby byla  $p(\beta | \tau, \mathbf{y})$  hustotou nedegenerovaného rozdělení, musí mít matice  $\mathbb{X}$  sloupcovou hodnot  $k$ .

Potom

$$\beta | \tau, \mathbf{y} \sim \mathcal{N}_k \left( \mathbf{b}(\mathbf{y}), \tau^{-1} (\mathbb{X}^\top \mathbb{X})^{-1} \right)$$



$$p(\tau | \mathbf{y}) \propto \cdots \propto \tau^{\frac{n-k}{2}-1} \exp\left(-\tau \frac{SS_e}{2}\right), \quad \tau > 0$$

- ▶ Aby existovala konečná normující konstanta, musí být  $\frac{n-k}{2} > 0$ , tj.  $n > k$ .  
Potom

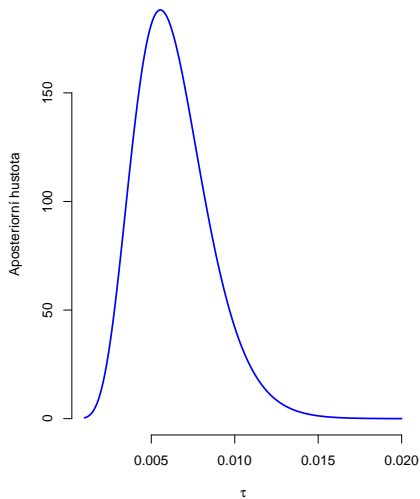
$$\tau | \mathbf{y} \sim \mathcal{G}\left(\frac{n-k}{2}, \frac{SS_e(\mathbf{y})}{2}\right)$$

- ▶ Rutinní použití věty o transformaci
  - ▣▶ aposteriorní hustota reziduálního rozptylu, resp. směrodatné odchylky

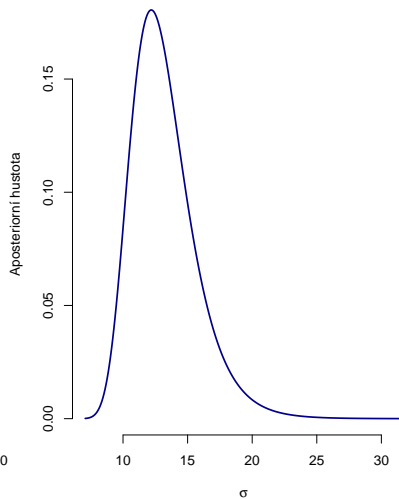
# Marginální aposteriorní rozdělení $\tau$ a $\sigma$

Vážení lehkých objektů,  $\tau | \mathbf{y} \sim \mathcal{G}(8, 1\,262.8)$

Reziduální přesnost



Reziduální směrodatná odchylka



- $\mathbb{E}(\tau | \mathbf{Y} = \mathbf{y}) = \frac{n-k}{SS_e}$
- $\text{var}(\tau | \mathbf{Y} = \mathbf{y}) = \frac{2(n-k)}{SS_e^2}$
- $\mathbb{E}(\sigma^2 | \mathbf{Y} = \mathbf{y}) = \frac{SS_e}{n-k-2}$ , pro  $n-k > 2$
- $\text{var}(\sigma^2 | \mathbf{Y} = \mathbf{y}) = \frac{2SS_e^2}{(n-k-2)^2(n-k-4)}$ , pro  $n-k > 4$
- $\mathbb{E}(\sigma | \mathbf{Y} = \mathbf{y}) = ???$
- $\text{var}(\sigma | \mathbf{Y} = \mathbf{y}) = ???$

## Bayesovské odhady $\tau$ a $\sigma$

---

Vážení lehkých objektů,  $\tau | \mathbf{y} \sim \mathcal{G}(8, 1\ 262.8)$

### (Možné) bodové odhady

$$\hat{\tau}_1 = \mathbb{E}(\tau | \mathbf{Y} = \mathbf{y}) = 0.00633 \quad \hat{\sigma}_1 = \mathbb{E}(\sigma | \mathbf{Y} = \mathbf{y}) = ???$$

$$\hat{\tau}_2 = \text{med}(\tau | \mathbf{Y} = \mathbf{y}) = 0.00607 \quad \hat{\sigma}_2 = \text{med}(\sigma | \mathbf{Y} = \mathbf{y}) = 12.83$$

### 95% věrohodnostní intervaly

$$\text{ET interval} \quad \tau : (0.00273, 0.01142) \quad \sigma : (9.36, 19.12)$$

$$\text{HPD interval} \quad \tau : (0.00235, 0.01079) \quad \sigma : (8.87, 18.22)$$

$\mathbf{T} \sim \text{MVT}_{k,\nu}(\boldsymbol{\Sigma})$ , jestliže

$$\mathbf{T} = \mathbf{U} \sqrt{\frac{\nu}{V}},$$

kde

- $\boldsymbol{\Sigma}$  je pozitivně definitní matice,

- $\mathbf{U} \sim \mathcal{N}_k(\mathbf{0}, \boldsymbol{\Sigma})$ ,

- $V \sim \chi_\nu^2$ ,

- $\mathbf{U}$  a  $V$  jsou nezávislé.

$\mathbf{T} \sim \text{MVT}_{k,\nu}(\mathbf{\Sigma})$  má hustotu

$$p(\mathbf{t}) = \frac{\Gamma\left(\frac{\nu+k}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right) \nu^{\frac{k}{2}} \pi^{\frac{k}{2}}} |\mathbf{\Sigma}|^{-\frac{1}{2}} \left\{1 + \frac{\mathbf{t}^\top \mathbf{\Sigma}^{-1} \mathbf{t}}{\nu}\right\}^{-\frac{\nu+k}{2}}, \quad \mathbf{t} \in \mathbb{R}^k$$

- $\mathbb{E}\mathbf{T} = \mathbf{0}$ , je-li  $\nu > 1$ ,
- $\text{var}\mathbf{T} = \frac{\nu}{\nu-2} \mathbf{\Sigma}$ , je-li  $\nu > 2$ ,
- $\text{modus}\mathbf{T} = \mathbf{0}$ .

Pro  $\boldsymbol{\mu} \in \mathbb{R}^k$ :

$\mathbf{Z} = \boldsymbol{\mu} + \mathbf{T}$ , kde  $\mathbf{T} \sim \text{MVT}_{k,\nu}(\boldsymbol{\Sigma})$  má hustotu

$$p(\mathbf{z}) = \frac{\Gamma(\frac{\nu+k}{2})}{\Gamma(\frac{\nu}{2}) \nu^{\frac{k}{2}} \pi^{\frac{k}{2}}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \left\{ 1 + \frac{(\mathbf{z} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu})}{\nu} \right\}^{-\frac{\nu+k}{2}}, \quad \mathbf{z} \in \mathbb{R}^k$$

- $\mathbb{E}\mathbf{Z} = \boldsymbol{\mu}$ , je-li  $\nu > 1$ ,
- $\text{var}\mathbf{Z} = \frac{\nu}{\nu-2}\boldsymbol{\Sigma}$ , je-li  $\nu > 2$ ,
- $\text{modus}\mathbf{Z} = \boldsymbol{\mu}$ .

$$p(\beta | \mathbf{y}) \propto \dots \propto \left\{ 1 + \frac{\mathbf{s}^{-2} (\beta - \mathbf{b})^\top \mathbb{X}^\top \mathbb{X} (\beta - \mathbf{b})}{n - k} \right\}^{-\frac{n-k+k}{2}}$$

- To jest,

$$\beta | \mathbf{y} \sim \mathbf{b}(\mathbf{y}) + \text{MVT}_{k, n-k} \left( \mathbf{s}^2(\mathbf{y}) (\mathbb{X}^\top \mathbb{X})^{-1} \right)$$

- Jaké je aposteriorní rozdělení pro  $\beta_j, j = 1, \dots, k$ ?

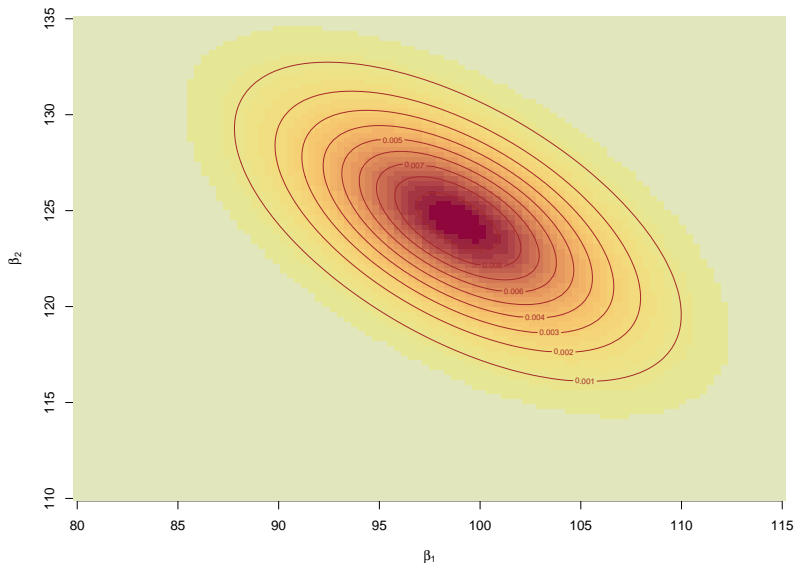
### Momenty

- $\mathbb{E}(\beta | \mathbf{Y} = \mathbf{y}) = \mathbf{b}, \quad \text{pro } n - k > 1$
- $\text{var}(\beta | \mathbf{Y} = \mathbf{y}) = \frac{n - k}{n - k - 2} \mathbf{s}^2 (\mathbb{X}^\top \mathbb{X})^{-1}, \quad \text{pro } n - k > 2$



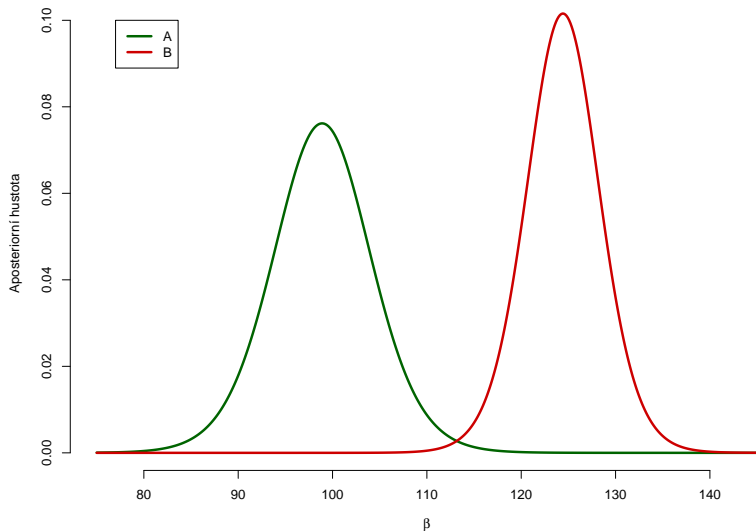
# Marginální aposteriorní rozdělení $\beta$

Vážení lehkých objektů



# Marginální aposteriorní rozdělení $\beta_1$ a $\beta_2$

Vážení lehkých objektů



Vážení lehkých objektů

(Možné) **bodové odhady**

$$\hat{\beta}_1 = \mathbb{E}(\beta_1 | \mathbf{Y} = \mathbf{y}) = \text{med}(\beta_1 | \mathbf{Y} = \mathbf{y}) = 98.89$$

$$\hat{\beta}_2 = \mathbb{E}(\beta_2 | \mathbf{Y} = \mathbf{y}) = \text{med}(\beta_2 | \mathbf{Y} = \mathbf{y}) = 124.42$$

**95% věrohodnostní intervaly**

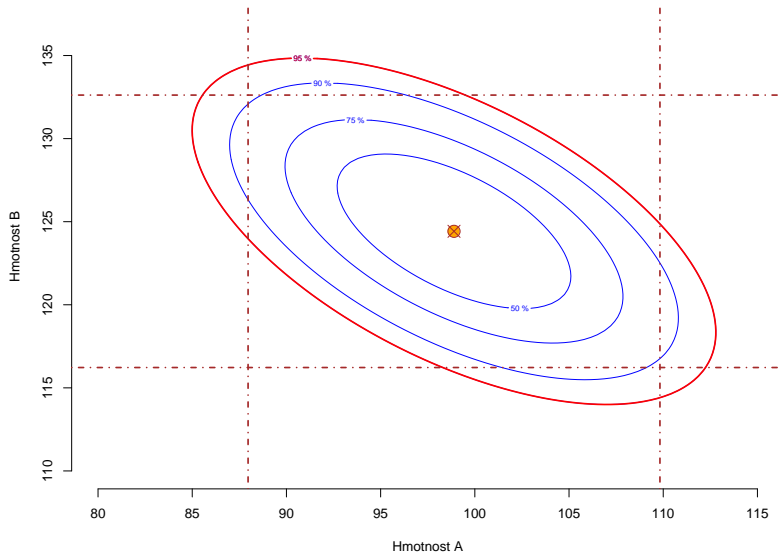
ET i HPD interval  $\beta_1 : (87.96, 109.83)$

$\beta_2 : (116.22, 132.62)$

$$\frac{s^{-2} (\beta - \mathbf{b})^T \mathbb{X}^T \mathbb{X} (\beta - \mathbf{b})}{k} \leq F_{k, n-k}(1 - \alpha)$$

# HPD věrohodnostní množina pro $\beta$

Vážení lehkých objektů



- $p(\beta, \tau) = p(\beta | \tau) \times p(\tau)$
- $p(\beta | \tau) \equiv \mathcal{N}_k(\beta_0, \tau^{-1} \Sigma_0)$
- $p(\tau) \equiv \mathcal{G}(c_0, d_0)$
- $\beta_0, \Sigma_0, c_0, d_0$  : pevné parametry apriorního rozdělání
  - Jaká je jejich interpretace?
  - Jakou apriorní informaci vyjadřují jejich konkrétní volby?

### Snadné domácí cvičení ☺:

- Odvoďte  $p(\beta, \tau | \mathbf{y})$ ,  $p(\beta | \mathbf{y})$ ,  $p(\tau | \mathbf{y})$  v případě, že je uvažováno výše uvedené apriorní rozdělání.

## Užitečná maticová rovnost

(pro snadné domácí cvičení)

Nechť  $\mathbb{A}$  a  $\mathbb{B}$  jsou symetrické pozitivně semidefinitní matice dimenze  $p \times p$  a alespoň jedna z nich necht' je pozitivně definitní. Necht'  $\mathbf{x}$ ,  $\mathbf{a}$ ,  $\mathbf{b}$  jsou reálné vektory z  $\mathbb{R}^p$ . Potom

$$\begin{aligned} & (\mathbf{x} - \mathbf{a})^\top \mathbb{A}(\mathbf{x} - \mathbf{a}) + (\mathbf{x} - \mathbf{b})^\top \mathbb{B}(\mathbf{x} - \mathbf{b}) \\ &= (\mathbf{x} - \boldsymbol{\mu})^\top \mathbb{V}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \\ & \quad \underbrace{- (\mathbb{A}\mathbf{a} + \mathbb{B}\mathbf{b})^\top \mathbb{V}(\mathbb{A}\mathbf{a} + \mathbb{B}\mathbf{b}) + \mathbf{a}^\top \mathbb{A}\mathbf{a} + \mathbf{b}^\top \mathbb{B}\mathbf{b}}_{\text{nezávisí na } \mathbf{x}}, \end{aligned}$$

kde

$$\begin{aligned} \mathbb{V} &= (\mathbb{A} + \mathbb{B})^{-1} \\ \boldsymbol{\mu} &= \mathbb{V}(\mathbb{A}\mathbf{a} + \mathbb{B}\mathbf{b}) \end{aligned}$$

# Konjugovaný systém apriorních rozdělání

## Aposteriorní rozdělání

$$p(\boldsymbol{\beta}, \tau | \mathbf{y}) \propto \tau^{\frac{n}{2} + \frac{k}{2} + c_0 - 1} \exp \left\{ -\tau \left( \frac{SS_e}{2} + d_0 \right) \right\} \\ \exp \left[ -\frac{\tau}{2} \left\{ (\boldsymbol{\beta} - \mathbf{b})^T \mathbb{X}^T \mathbb{X} (\boldsymbol{\beta} - \mathbf{b}) + (\boldsymbol{\beta} - \boldsymbol{\beta}_0)^T \boldsymbol{\Sigma}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\beta}_0) \right\} \right],$$

$$\tau | \mathbf{y} \sim \mathcal{G} \left( \frac{n}{2} + c_0, \frac{SS_e(\mathbf{y})}{2} + d_0 \right),$$

$$\boldsymbol{\beta} | \tau, \mathbf{y} \sim \mathcal{N}_k(\hat{\boldsymbol{\beta}}(\mathbf{y}), \tau^{-1} \hat{\boldsymbol{\Sigma}}(\mathbf{y})), \quad \hat{\boldsymbol{\Sigma}}(\mathbf{y}) = (\mathbb{X}^T \mathbb{X} + \boldsymbol{\Sigma}_0^{-1})^{-1}, \\ \hat{\boldsymbol{\beta}}(\mathbf{y}) = \hat{\boldsymbol{\Sigma}}(\mathbf{y}) (\mathbb{X}^T \mathbb{X} \mathbf{b}(\mathbf{y}) + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\beta}_0)$$

Interpretace parametrů apriorního rozdělání?



# 2

**Aposteriorní inference založená na simulaci**

## Oddíl 2.1

# Úvod

# Obecná formulace problému

---

- $\theta = (\theta_1, \dots, \theta_k)^\top$  : vektor (bayesovských) parametrů
- Aposteriorní inference je založená na

$$p(\theta | \mathbf{y}) \propto L(\theta) p(\theta)$$

- hustota vzhledem k  $\sigma$ -konečné míře  $\lambda$  na  $(\mathbb{R}^k, \mathcal{B}^k)$ ,
- v dalším budeme zkráceně namísto  $p(\theta | \mathbf{y})d\lambda(\theta)$  psát  $p(d\theta | \mathbf{y})$ .

- V rámci posteriorní inference nás též pro nějaké měřitelné funkce  $t$  na parametrickém prostoru  $\Theta$  zajímají zejména následující veličiny:
  - $\mathbb{E}_{p(d\theta | \mathbf{y})} t(\theta) = \int t(\theta) p(d\theta | \mathbf{y})$
  - kvantily  $t(\theta)$  vzhledem k rozdělení  $p(d\theta | \mathbf{y})$ 
    - ▮ medián, věrohodnostní intervaly,
  - posteriorní rozdělení náhodné veličiny.  $t(\theta)$
- Kromě triviálních případů vyžaduje **integrování** přes  $p(d\theta | \mathbf{y})$ .

## Oddíl 2.2

# Monte Carlo integrace

Snaha pro měřitelnou funkci  $t : \mathbb{R}^k \rightarrow \mathbb{R}$  **numericky** spočítat

$$\mathbb{E}_{p(d\theta | \mathbf{y})} t(\theta) = \int t(\theta) p(d\theta | \mathbf{y}).$$

- Předpokládejme, že  $\mathbb{E}_{p(d\theta | \mathbf{y})} t(\theta)$  existuje konečné.
- Předpokládejme dále, že jsme schopni získat **náhodný výběr**

$$S_M = \{\theta^{(1)}, \dots, \theta^{(M)}\} \quad \text{z rozdělení } p(d\theta | \mathbf{y}).$$

$$\text{Potom } \int t(\theta) p(d\theta | \mathbf{y}) \approx \frac{1}{M} \sum_{m=1}^M t(\theta^{(m)}) = \widehat{\mathbb{E}}_{p(d\theta | \mathbf{y})} t(\theta) := \widehat{t}_M.$$

- Proč a kdy to funguje?

- Jestliže  $\int |t(\theta)| p(d\theta | \mathbf{y}) < \infty$ , potom

$$\hat{t}_M = \frac{1}{M} \sum_{m=1}^M t(\theta^{(m)}) \xrightarrow{\text{a.s.}} \int t(\theta) p(d\theta | \mathbf{y}) \quad \text{pro } M \rightarrow \infty$$

(silný zákon velkých čísel).

- Dále, jestliže  $\int \{t(\theta)\}^2 p(d\theta | \mathbf{y}) < \infty$ , potom

$$\begin{aligned} v_M := \text{var} \left\{ \frac{1}{M} \sum_{m=1}^M t(\theta^{(m)}) \right\} &= \frac{1}{M} \text{var}_{p(d\theta | \mathbf{y})} t(\theta) \\ &= \frac{1}{M} \int \{t(\theta) - \mathbb{E}_{p(\theta | \mathbf{y})}\}^2 p(d\theta | \mathbf{y}) \end{aligned}$$

ohodnocuje rychlost konvergence a tudíž též přesnost aproximace při použití konečného  $M$ .

- $\sqrt{v_M}$  = Monte Carlo chyba (*Monte Carlo Error*).

# Monte Carlo integrace

## Odhad Monte Carlo chyby

- Opět s využitím **silného zákona velkých čísel** se snadno ukáže (obdobně jako při důkazu konzistence výběrového rozptylu), že

$$\frac{1}{M-1} \sum_{m=1}^M \left\{ t(\theta^{(m)}) - \hat{t}_M \right\}^2 \xrightarrow{\text{a.s.}} \text{var}_{p(d\theta | y)} t(\theta), \quad \text{pro } M \rightarrow \infty.$$

- Rozptyl  $v_M$  lze tudíž odhadnout pomocí

$$\hat{v}_M = \frac{1}{M(M-1)} \sum_{m=1}^M \left\{ t(\theta^{(m)}) - \hat{t}_M \right\}^2.$$

- $\sqrt{\hat{v}_M}$  = odhad Monte Carlo chyby.



- Za položených předpokladů platí též **centrální limitní věta**, která spolu s **Cramérovou-Sluckého větou** vede k

$$\frac{\hat{t}_M - \mathbb{E}_{p(d\theta | \mathbf{y})} t(\boldsymbol{\theta})}{\sqrt{\hat{v}_M}} \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1)$$

- ohodnocení konvergence Monte Carlo integrace,
- konstrukce konfidenčních mezí pro aproximaci.

### Poznámka

- Potřebný náhodný výběr  $\mathcal{S}_M$  obvykle na náš popud generuje počítač.
- Velikost výběru  $M$  lze poměrně snadno ovlivnit (více méně závisí pouze na výkonu počítače a době, po kterou jsme ochotni čekat na výsledek).
- S patřičně velkým  $M$  lze i reálně dosáhnout požadované přesnosti při odhadu  $\mathbb{E}_{p(d\theta | \mathbf{y})} t(\boldsymbol{\theta})$ .

## Oddíl 2.3

# Důležité speciální případy

## Důležité speciální případy

- Předpokládejme, že  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)^\top$  a  $\mathcal{S}_M = \{\boldsymbol{\theta}^{(1)}, \dots, \boldsymbol{\theta}^{(M)}\}$  je náhodný výběr z  $p(d\boldsymbol{\theta} | \mathbf{y})$

### Důležitý speciální případ 1

$$t(\boldsymbol{\theta}) = \theta_j \quad (j = 1, \dots, k)$$

- $\frac{1}{M} \sum_{m=1}^M \theta_j^{(m)}$  je MC odhadem pro  $\mathbb{E}_{p(d\boldsymbol{\theta} | \mathbf{y})} \theta_j$ .

## Důležitý speciální případ 2

$$t(\theta) = \mathbb{I}_{[\theta_j \leq x]}(\theta_j) \text{ pro } x \in \mathbb{R} \quad (j = 1, \dots, k)$$

- $\frac{1}{M} \sum_{m=1}^M \mathbb{I}_{[\theta_j \leq x]}(\theta_j^{(m)}) = \frac{\#[\theta_j^{(m)} \leq x]}{M}$  je MC odhadem pro  $\mathbb{E}_{p(d\theta | \mathbf{y})} \mathbb{I}_{[\theta_j \leq x]}(\theta_j) = \mathbf{P}(\theta_j \leq x | \mathbf{Y} = \mathbf{y})$ .
- MC odhad **marginální** aposteriorní distribuční funkce  $\theta_j$ .
- Libovolný odhad hustoty (histogram, jádrový odhad, ...) založený na  $\theta_j^{(1)}, \dots, \theta_j^{(M)}$  je MC odhadem **marginální** aposteriorní hustoty  $p(\theta_j | \mathbf{y})$ ,  
▮▮▮ odsud lze mj. odhadnout **HPD věrohodnostní intervaly**.
- Výběrové kvantily založené na  $\theta_j^{(1)}, \dots, \theta_j^{(M)}$  odhadují kvantily **marginálního** aposteriorního rozdělení  $p(d\theta_j | \mathbf{y})$   
▮▮▮ odhady **ET věrohodnostních intervalů**.

## Důležitý speciální případ 3

Veličina hlavního zájmu je  $r(\theta)$ , kde  $r : \mathbb{R}^k \rightarrow \mathbb{R}$  je měřitelná funkce  
a  $t(\theta) = \mathbb{I}_{[r(\theta) \leq x]}(\theta)$  pro  $x \in \mathbb{R}$

- $\frac{1}{M} \sum_{m=1}^M \mathbb{I}_{[r(\theta) \leq x]}(\theta^{(m)}) = \frac{\#[r(\theta^{(m)}) \leq x]}{M}$  je MC odhadem pro

$$\mathbb{E}_{p(d\theta | \mathbf{y})} \mathbb{I}_{[r(\theta) \leq x]}(\theta) = P\{r(\theta) \leq x \mid \mathbf{Y} = \mathbf{y}\}.$$

- MC odhad **marginální** aposteriorní distribuční funkce  $r(\theta)$ .
- Libovolný odhad hustoty (histogram, jádrový odhad, ...) založený na  $r(\theta^{(1)}), \dots, r(\theta^{(M)})$  je MC odhadem **marginální** aposteriorní hustoty  $p(r(\theta) \mid \mathbf{y})$ 
  - odsud lze mj. odhadnout **HPD věrohodnostní intervaly**.
- Výběrové kvantily založené na  $r(\theta^{(1)}), \dots, r(\theta^{(M)})$  odhadují kvantily **marginálního** aposteriorního rozdělení  $p(dr(\theta) \mid \mathbf{y})$ 
  - odhady **ET věrohodnostních intervalů**.

### Důležitý speciální případ 3, pokračování

Veličina hlavního zájmu je  $r(\theta)$ , kde  $r : \mathbb{R}^k \rightarrow \mathbb{R}$  je měřitelná funkce

- Zvolíme-li  $t(\theta) = r(\theta)$ , dostaneme MC odhad pro  $\mathbb{E}_{p(d\theta | \mathbf{y})} r(\theta)$ .
- Při počítání MC odhadů veličin založených na aposteriorním rozdělení  $r(\theta)$  již nemusíme ani integrovat, ani používat větu o transformaci.
- $r(\theta)$  může být prakticky libovolně “složitou” funkcí  $\theta$  a stále je vše upočítatelné i prakticky.

Kde vzít onen náhodný výběr

$$S_M = \{\boldsymbol{\theta}^{(1)}, \dots, \boldsymbol{\theta}^{(M)}\} \quad \text{z rozdělení } p(d\boldsymbol{\theta} | \mathbf{y})?$$

## Oddíl 2.4

# Simulace ze zadaného rozdělení



- Vše dosud řečené o MC integraci lze samozřejmě použít k numerickému počítání integrálů vzhledem k **libovolnému rozdělení**.
  - Nejen aposterioriálnímu v rámci nějakého bayesovského modelu.
- V dalším se budeme zabývat problémem simulace ze zadaného rozdělení, které budeme reprezentovat hustotou  $f(\theta)$  vzhledem k nějaké  $\sigma$ -konečné míře  $\lambda$  na  $(\mathbb{R}^k, \mathcal{B}^k)$  a distribuční funkcí  $F(\theta)$ .
  - Nadále budeme zkráceně psát  $f(\theta)d\lambda(\theta) = f(d\theta)$ .
  - V rámci bayesovských metod budeme obvykle používat s  $f(d\theta) = p(d\theta | \mathbf{y})$ .

- Nejprve se budeme zabývat případem **jednorozměrného** rozdělení  $f(d\theta)$ .

### Věta 2.1 .

---

*Nechť  $U$  je náhodná veličina s rovnoměrným rozdělením na intervalu  $(0, 1)$ .  
Nechť  $\theta$  je náhodná veličina s distribuční funkcí  $F$ . Nechť pro  $0 < u < 1$  je  $F^{-1}(u) = \inf\{\theta : F(\theta) \geq u\}$  kvantilová funkce. Potom má náhodná veličina*

$$Z = F^{-1}(U)$$

*stejné rozdělení jako náhodná veličina  $\theta$ , tj.  $Z$  má rozdělení s distribuční funkcí  $F$ .*

*Důkaz.*

- Triviální, jestliže navíc předpokládáme, že  $F$  je **spojitá** a **rostoucí** na nosiči, tj. na množině  $\Theta$ , pro kterou platí

$$P(\theta \in \Theta) = 1, \quad \forall \tilde{\Theta} \subset \Theta, P(\theta \in \Theta \setminus \tilde{\Theta}) > 0 \Rightarrow P(\theta \in \tilde{\Theta}) < 1.$$

- O něco složitější, jestliže spojitost a monotonii  $F$  na nosiči nepředpokládáme.



# Univerzální algoritmus pro simulování z jednorozměrného rozdělení

---

Generování z rozdělení  $\theta$  s distribuční funkcí  $F(\theta)$ :

1. Vygeneruj  $U \sim \mathcal{U}(0, 1)$ .
  2. Polož  $\theta = F^{-1}(U)$ .
- Efektivní, jestliže jsme schopni efektivně počítat hodnoty  $F^{-1}$ .

## Specifické algoritmy pro simulování z jednorozměrných rozdělání

- Pro mnohá běžná rozdělání existují efektivnější metody založené obvykle na transformacích, z nichž mnohé byly propočítány v rámci **Cvičení k předmětu Matematická statistika 1 (NMSA331)**.

- Příklad 1:

$U_1, U_2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}(0, 1)$ , potom pro

$$\theta_1 = \mu + \sigma \cos(2\pi U_1) \sqrt{-2 \log U_2},$$

$$\theta_2 = \mu + \sigma \sin(2\pi U_1) \sqrt{-2 \log U_2}$$

platí  $\theta_1, \theta_2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, \sigma^2)$ .

- Příklad 2:

$U \sim \mathcal{U}(0, 1)$ , potom pro

$$\theta = -\frac{1}{\lambda} \log(U)$$

platí  $\theta \sim \text{Exp}(\lambda)$ .


## Více o simulování z rozdělení s plně specifikovanou hustotou (nejenom jednorozměrnou)

---

### Podrobnější obrázek

- Přednáška **Simulační metody (NMST535)** vyučovaná v zimním semestru.
- Luc Devroye. *Non-Uniform Random Variate Generation*. New York: Springer-Verlag, 1986.

### Praktické aplikace

- Statistické balíky využívají těchto metod ke generování z většiny běžných rozdělení.
  -  funkce `runif`, `rnorm`, `rexp`, ...
  - “Běžný” uživatel se nemusí příliš trápit s generováním z běžných ( $\equiv$  “pojmenovaných”) rozdělení a prostě použije zabudované možnosti svého software.

- Předpokládejme, že  $\theta = (\theta_1^\top, \dots, \theta_k^\top)^\top$  a v rozkladu

$$\begin{aligned} f(d\theta) &= f(d\theta_1, \dots, d\theta_k) \\ &= f(d\theta_1 \mid \theta_2, \dots, \theta_k) f(d\theta_2 \mid \theta_3, \dots, \theta_k) \cdots f(d\theta_{k-1} \mid \theta_k) f(d\theta_k) \end{aligned}$$

jsme schopni snadno generovat ze všech (podmíněných) rozdělení  $f(d\theta_1 \mid \theta_2, \dots, \theta_k)$ ,  $f(d\theta_2 \mid \theta_3, \dots, \theta_k)$ ,  $\dots$ ,  $f(d\theta_{k-1} \mid \theta_k)$ ,  $f(d\theta_k)$ .

- Algoritmus pro generování ze sdruženého rozdělení  $f(d\theta)$  je potom následující:
  1. Vygeneruj  $\theta_k$  z  $f(d\theta_k)$ .
  2. Vygeneruj  $\theta_{k-1}$  z  $f(d\theta_{k-1} \mid \theta_k)$ .
  3.  $\vdots$
  4. Vygeneruj  $\theta_2$  z  $f(d\theta_2 \mid \theta_3, \dots, \theta_k)$ .
  5. Vygeneruj  $\theta_1$  z  $f(d\theta_1 \mid \theta_2, \theta_3, \dots, \theta_k)$ .

## Lineární model

$$\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \mathbb{X} : \text{pevná matice } n \times k,$$
$$\boldsymbol{\varepsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n).$$

- Parametry:  $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \tau^\top)^\top$ , kde  $\tau = \sigma^{-2} > 0$ .
- Věrohodnost:  $\mathbf{Y} \sim \mathcal{N}_n(\mathbb{X}\boldsymbol{\beta}, \tau^{-1} \mathbf{I}_n)$ .
- Neinformativní apriorní rozdělení:

$$p(\boldsymbol{\beta}) \propto 1, \quad \boldsymbol{\beta} \in \mathbb{R}^k,$$
$$p(\tau) \propto \frac{1}{\tau}, \quad \tau > 0.$$



## Příklad: Lineární model s neinformativní apriorním rozdělením

### Aposteriorní rozdělení

- Označme:  $\mathbf{b} = (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbf{y}$ ,  
 $SS_e = \|\mathbf{y} - \mathbb{X}\mathbf{b}\|^2$ .

- Bylo odvozeno:

$$p(\boldsymbol{\beta}, \tau | \mathbf{y}) = p(\boldsymbol{\beta} | \tau, \mathbf{y}) \times p(\tau | \mathbf{y}),$$

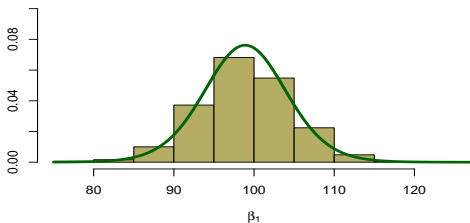
kde  $p(\tau | \mathbf{y}) \sim \mathcal{G}\left(\frac{n-k}{2}, \frac{SS_e}{2}\right)$ ,

$$p(\boldsymbol{\beta} | \tau, \mathbf{y}) \sim \mathcal{N}_k\left(\mathbf{b}, \tau^{-1} (\mathbb{X}^\top \mathbb{X})^{-1}\right).$$

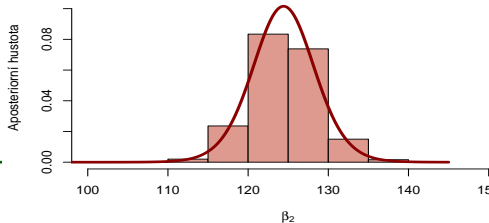
# Příklad: Vážení lehkých objektů

Marginální aposteriorní hustoty ( $M=1\ 000$ )

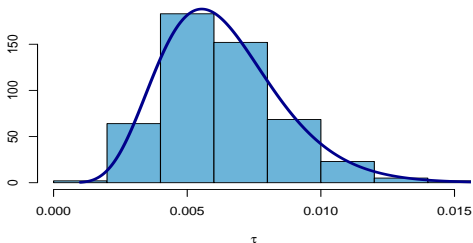
Hmotnost A



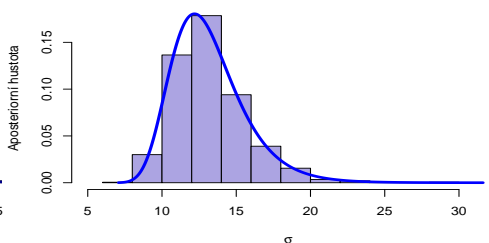
Hmotnost B



Inverzní rozptyl



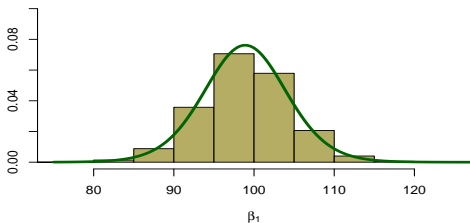
Sm rodatná odchylka



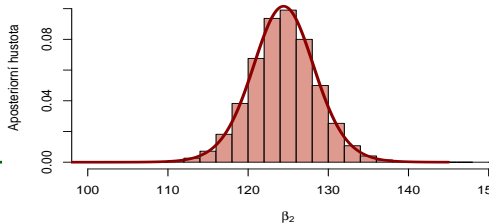
# Příklad: Vážení lehkých objektů

Marginální aposteriorní hustoty ( $M=1\ 000\ 000$ )

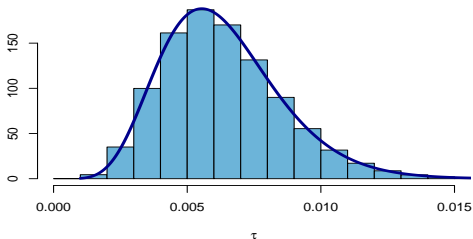
Hmotnost A



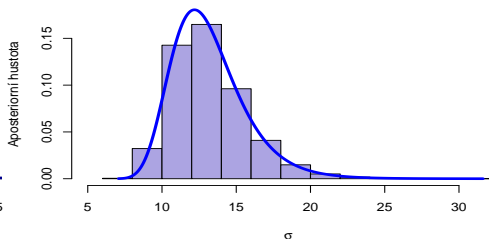
Hmotnost B



Inverzní rozptyl



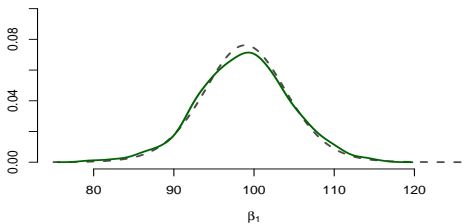
Sm rodatná odchylka



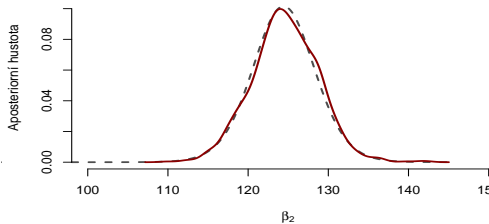
# Příklad: Vážení lehkých objektů

Marginální aposteriorní hustoty (M=1 000)

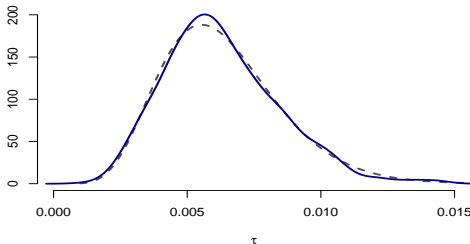
Hmotnost A



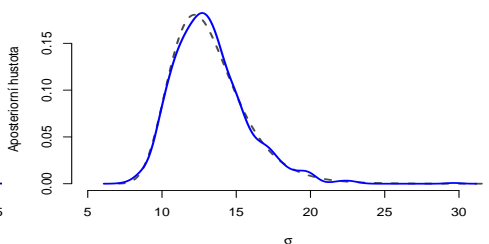
Hmotnost B



Inverzní rozptyl



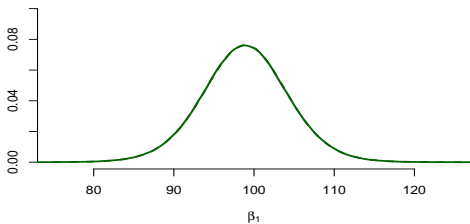
Inverzní rozptyl



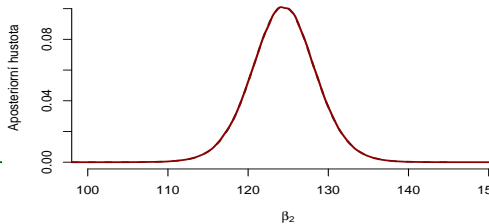
# Příklad: Vážení lehkých objektů

Marginální aposteriorní hustoty ( $M=1\ 000\ 000$ )

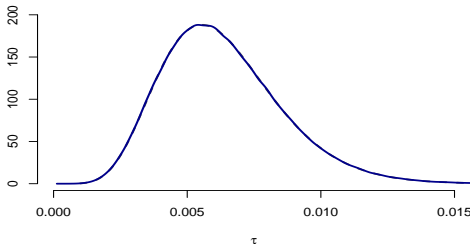
Hmotnost A



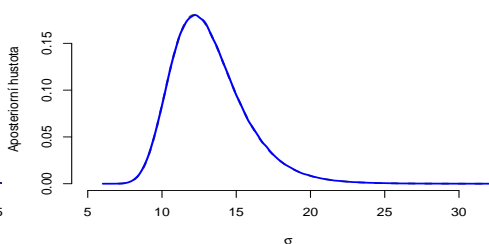
Hmotnost B



Inverzní rozptyl

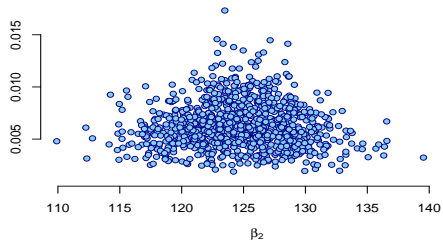
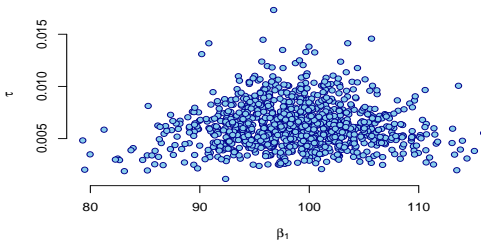
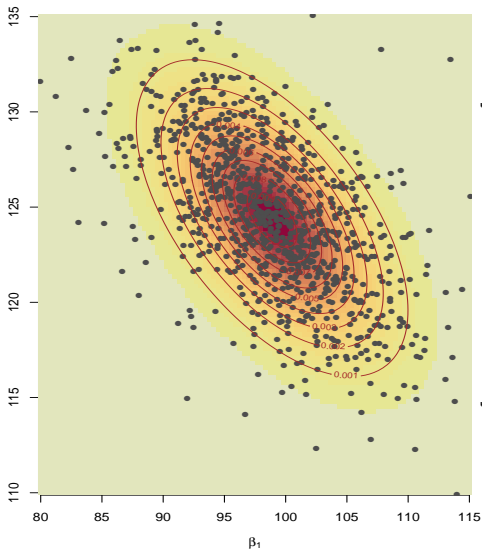


Inverzní rozptyl



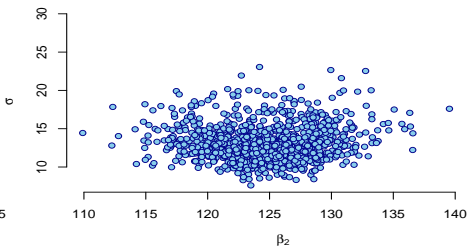
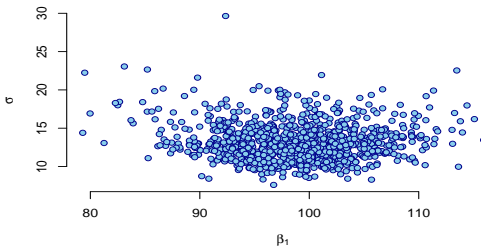
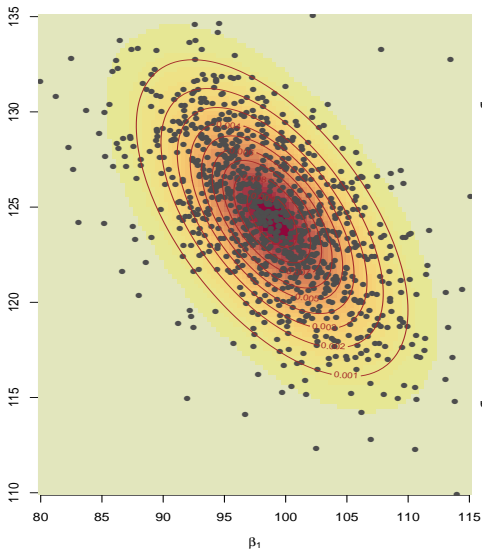
# Příklad: Vážení lehkých objektů

Sdružené výběry z aposteriorního rozdělení ( $M=1\ 000$ )



# Příklad: Vážení lehkých objektů

Sdružené výběry z aposteriorního rozdělení ( $M=1\ 000$ )



## Příklad: Vážení lehkých objektů

Apsteriorní inference pro  $\beta$  (M=1 000)

	$\beta_1$	$\beta_2$
<b>Aposter. střední hodnota</b>	98,8947	124,4211
<b>MC odhad</b>	98,7944	124,6197
<b>MC chyba</b>	0,1804	0,1312
<b>Aposter. medián</b>	98,8947	124,4211
<b>MC odhad</b>	98,7813	124,5673
<b>95% ET věr. interval</b>	(87,9641; 109,8253)	(116,2231; 132,6190)
<b>MC odhad</b>	(86,9761; 110,0815)	(116,7594; 132,6802)
<b>95% HPD věr. interval</b>	(87,9641; 109,8253)	(116,2231; 132,6190)
<b>MC odhad</b>	(87,9245; 110,7076)	(116,4892; 132,2210)



## Příklad: Vážení lehkých objektů

Apsteriorní inference pro  $\beta$  ( $M=1\ 000\ 000$ )

	$\beta_1$	$\beta_2$
<b>Apster. střední hodnota</b>	98,8947	124,4211
<b>MC odhad</b>	98,8874	124,4192
<b>MC chyba</b>	0,0055	0,0041
<b>Apster. medián</b>	98,8947	124,4211
<b>MC odhad</b>	98,8849	124,4196
<b>95% ET věr. interval</b>	(87,9641; 109,8253)	(116,2231; 132,6190)
<b>MC odhad</b>	(87,9680; 109,8184)	(116,2239; 132,6113)
<b>95% HPD věr. interval</b>	(87,9641; 109,8253)	(116,2231; 132,6190)
<b>MC odhad</b>	(88,0765; 109,9202)	(116,1496; 132,5329)

## Příklad: Vážení lehkých objektů

Apsteriorní inference pro  $\tau$  a  $\sigma$  ( $M=1\ 000$ )

	$\tau$	$\sigma$
<b>Apster. střední hodnota</b>	0,00633	?
<b>MC odhad</b>	0,00629	13,207
<b>MC chyba</b>	0,0000689	0,0776
<b>Apster. medián</b>	0,00607	?
<b>MC odhad</b>	0,00601	12,904
<b>95% ET věr. interval</b>	(0,00273; 0,01142)	?
<b>MC odhad</b>	(0,00272; 0,01097)	(9,547; 19,183)
<b>95% HPD věr. interval</b>	?	?
<b>MC odhad</b>	(0,00248; 0,01039)	(8,987; 18,186)

## Příklad: Vážení lehkých objektů

Apsteriorní inference pro  $\tau$  a  $\sigma$  ( $M=1\,000\,000$ )

	$\tau$	$\sigma$
<b>Apster. střední hodnota</b>	0,00633	?
<b>MC odhad</b>	0,00633	13,198
<b>MC chyba</b>	0,0000022	0,0025
<b>Apster. medián</b>	0,00607	?
<b>MC odhad</b>	0,00607	12,838
<b>95% ET věr. interval</b>	(0,00273; 0,01142)	?
<b>MC odhad</b>	(0,00274; 0,01142)	(9,356; 19,116)
<b>95% HPD věr. interval</b>	?	?
<b>MC odhad</b>	(0,00237; 0,01081)	(8,872; 18,224)

# 3

## **Hierarchické modely**

## Oddíl **3.1**

# Hierarchické apriorní rozdělení

- Volba apriorního rozdělení může značným způsobem ovlivnit formu rozdělení **aposteriorního**.
- Nebezpečí **zneužití** bayesovské statistiky.
- Většina “bayesovských” aplikací z posledních cca 25 let
  - není motivována snahou využívat jakoukoliv apriorní informaci,
  - hlavní motivace: navržený model nelze (ani numericky) odhadnout frekventisticky (typicky pomocí **maximální věrohodnosti**), nicméně lze ho odhadnout pomocí **simulací** bayesovsky,
  - neexistuje žádná skutečná apriorní informace.

- Pouze málokdy je apriorní informace dostatečně bohatá na to, abychom mohli zvolené apriorní rozdělení považovat za přesně a **bez jakékoliv chyby** definované.
- Potřeba vhodným způsobem vyjádřit **nejistotu** při volbě apriorního rozdělení.
- ▶ **Bayesovský model s hierarchicky specifikovaným apriorním rozdělením**
  - rozklad apriorního rozdělení do několika úrovní podmíněných rozdělení,
  - nejistota na libovolné úrovni je vyjádřena apriorním rozdělením v další úrovni.

# Bayesovský model s hierarchicky specifikovaným apriorním rozdělením

**Definice 3.1** Bayesovský model s hierarchicky specifikovaným apriorním rozdělením.

Bayesovský model s hierarchicky specifikovaným apriorním rozdělením je statistický model s věrohodností  $L(\psi) = p(\mathbf{y} | \psi)$  a apriorním rozdělením  $p(\psi)$ , kde  $p(\psi)$  je rozloženo na podmíněná rozdělení

$$p_0(\psi | \zeta_1), p_1(\zeta_1 | \zeta_2), \dots, p_{m-1}(\zeta_{m-1} | \zeta_m)$$

a marginální rozdělení  $p_m(\zeta_m)$  tak, že

$$p(\psi) =$$

$$\int_{Z_1 \times \dots \times Z_m} p_0(\psi | \zeta_1) p_1(\zeta_1 | \zeta_2) \cdots p_{m-1}(\zeta_{m-1} | \zeta_m) p_m(\zeta_m) d\zeta_1 \cdots d\zeta_m.$$

Parametry obsažené v  $\zeta_i$  se nazývají **hyperparametry**  $i$ -té úrovně ( $1 \leq i \leq m$ ).



### Model

$$\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \mathbb{X} : \text{pevná matice } n \times k,$$
$$\boldsymbol{\varepsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n).$$

- **Parametry:**  $\boldsymbol{\psi} = (\boldsymbol{\beta}^\top, \tau)^\top$ , kde  $\tau = \sigma^{-2} > 0$
- **Věrohodnost:**  $L(\boldsymbol{\psi}) = p(\mathbf{y} | \boldsymbol{\psi}) \equiv \mathcal{N}_n(\mathbb{X}\boldsymbol{\beta}, \tau^{-1} \mathbf{I}_n)$
- **Konjugované apriorní rozdělení:**

$$p(\boldsymbol{\beta}, \tau) = p(\boldsymbol{\beta} | \tau) \times p(\tau)$$

$$p(\boldsymbol{\beta} | \tau) \equiv \mathcal{N}_k(\boldsymbol{\beta}_0, \tau^{-1} \boldsymbol{\Sigma}_0)$$

$$p(\tau) \equiv \mathcal{G}(c_0, d_0)$$

- $\boldsymbol{\beta}_0, \boldsymbol{\Sigma}_0, c_0, d_0$  : pevné (hyper)parametry.

## Příklad: Lineární model

- Volba parametrů gama apriorního rozdělení pro  $\tau$  (resp. samotná volba gama rozdělení) mívá dosti značný vliv na výsledné aposteriorní rozdělení.
- Nepovažujme  $c_0$  a/nebo  $d_0$  za pevné konstanty, ale umožněme náhodnost při jejich výběru.

### ▣ hierarchický model

- například:

$$p(\tau | d_0) \equiv \mathcal{G}(c_0, d_0)$$

$$p(d_0) \equiv \mathcal{G}(g_0, h_0)$$

- $c_0$  : pevný (hyper)parametr
- $d_0$  : náhodný hyperparametr 1. úrovně
- $g_0, h_0$  : pevné (hyper)parametry (2. úrovně)

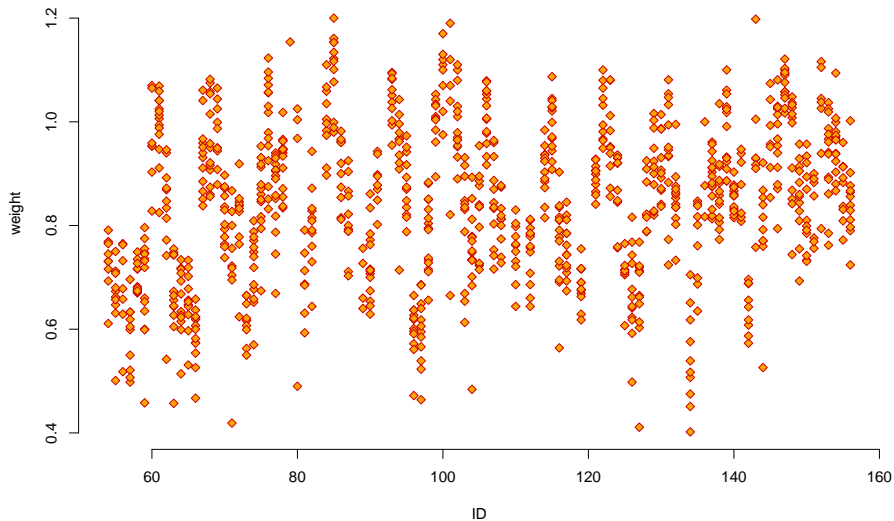
## Oddíl **3.2**

# Hierarchicky specifikovaná věrohodnost

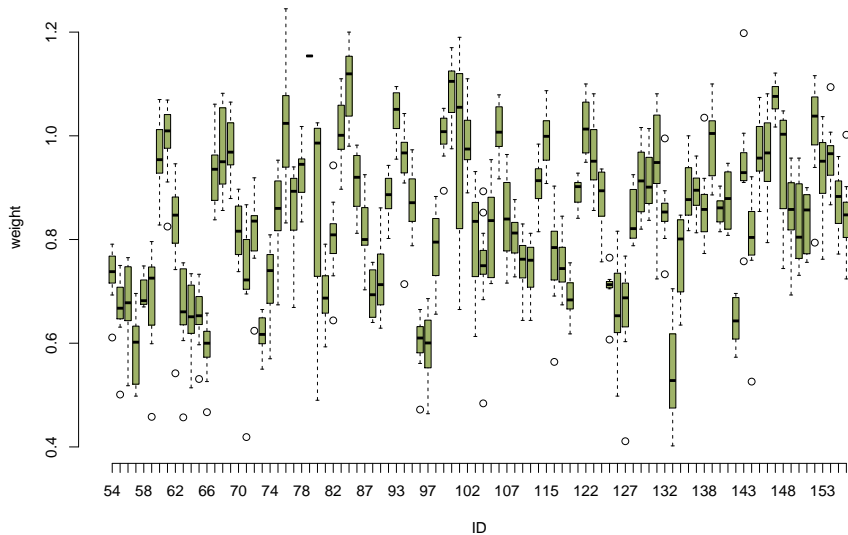
- Hierarchická specifikace (též v nebayesovském kontextu) je často přirozeným způsobem jak zkonstruovat **realistický** pravděpodobnostní model pro popis reálné situace.
- Oblasti využití, které jste už potkali
  - Data získaná **stratifikovaným** výběrem či jinak **shlukovaná** (*grouped data*).
  - **Longitudinální** (biostatistika), resp. **panelová** (ekonometrie) data.

- Data z *National Toxicology Program*
- 94 těhotným myším byla podána v předem určených momentech určená množství etylenglykolu (EG)
- V 17. dnu těhotenství byly myši usmrceny a následně byla zaznamenána hmotnost zárodků
- **Pravděpodobnostní reprezentace dat:**  
 $Y_{i,j}$  = hmotnost  $j$ -tého zárodku  $i$ -té myši,  $i = 1, \dots, N$ ,  $j = 1, \dots, n_i$
- **Primární cíl:**  
Odhad a inference pro  $\mu = \mathbb{E}Y_{i,j}$

## Příklad: NTP TER84073 pokus na myších



# Příklad: NTP TER84073 pokus na myších



- Možný pravděpodobnostní model:

$$Y_{i,j} \sim \mathcal{N}(\mu, \sigma^2).$$

- 
- Je ospravedlnitelné předpokládat, že  $Y_{1,1}, \dots, Y_{N,n_N}$  jsou nezávislé?



## Příklad: NTP TER84073 pokus na myších

- Realističtější pravděpodobnostní model:

$$Y_{i,j} | b_i \sim \mathcal{N}(b_i, \sigma^2),$$
$$b_i \sim \mathcal{N}(\mu, d^2).$$

- Pro každé  $i$  již lze předpokládat (podmíněnou) nezávislost  $Y_{i,1} | b_i, \dots, Y_{i,n_i} | b_i$ .
- Lze též předpokládat nezávislost  $b_1, \dots, b_N$ .
- $b_i$  : střední hmotnost zárodku  $i$ -té myši.
- $\sigma^2$  : rozptyl hmotnosti zárodků u jednotlivé myši  
▮ vnitroskupinová variabilita.
- $\mu$  : střední hmotnost zárodku v celé populaci

$$\mathbb{E}Y_{i,j} = \mathbb{E}\{\mathbb{E}(Y_{i,j} | b_i)\} = \mathbb{E}b_i = \mu.$$

## Příklad: NTP TER84073 pokus na myších

---

- Realističtější pravděpodobnostní model:

$$Y_{i,j} | b_i \sim \mathcal{N}(b_i, \sigma^2),$$
$$b_i \sim \mathcal{N}(\mu, d^2).$$

- Modelujeme též jistým způsobem korelaci mezi zárodky jedné myši:

$$\text{cov}(Y_{i,j}, Y_{i,k}) = \dots = d^2,$$
$$\text{var}(Y_{i,j}) = \dots = \sigma^2 + d^2.$$

$$\Rightarrow \text{cor}(Y_{i,j}, Y_{i,k}) = \frac{d^2}{\sigma^2 + d^2}$$

⇒ **vnitroskupinová korelace** (*intraclass correlation*)



# Hierarchické modely

---

## Obecná poznámka

- Též v jiných modelech lze často rozlišit tři typy parametrů (v bayesovském smyslu):
  - **skrytá data**
    - ▮ v dalším budeme obvykle značit  $\xi$
  - **“čisté” parametry**  $\equiv$  parametry též ve frekventistickém pojetí
    - ▮ v dalším budeme obvykle značit  $\psi$
  - **náhodné hyperparametry**
    - ▮ v dalším budeme obvykle značit  $\zeta$
- Parametry pro **bayesovský** model jsou potom  $\theta = (\xi^\top, \psi^\top, \zeta^\top)^\top$ .

# Hierarchické modely

---

## Obecná poznámka, pokračování

- **Sdružené** apriorní rozdělení je zadáno rozkladem

$$p(\theta) = p(\xi, \psi, \zeta) = p(\xi | \psi) p(\psi | \zeta) p(\zeta)$$

- $p(\xi | \psi)$  : strukturální část apriorního rozdělení
  - ▮ plyne z uvažovaného pravděpodobnostního modelu použitého pro popis situace
- $p(\psi | \zeta) p(\zeta)$  : “skutečné” apriorní rozdělení

## “ANOVA” hierarchický model

Konkrétní apriorní rozdělení (jedna z možností)

- Např. **konjugované apriorní** rozdělení s dodatečnými hyperparametry, abychom se ochránili od nadměrného ovlivnění aposteriorního rozdělení rozdělením apriorním:

$$p(\tau, \mu, q, b_0, d_0) = \underbrace{p(\mu | q)}_{\mathcal{N}(\mu_0, k_0^{-1} q^{-1})} \underbrace{p(q | b_0)}_{\mathcal{G}(a_0, b_0)} \underbrace{p(b_0)}_{\mathcal{G}(p_0, r_0)} \underbrace{p(\tau | d_0)}_{\mathcal{G}(c_0, d_0)} \underbrace{p(d_0)}_{\mathcal{G}(g_0, h_0)}$$

- $b_0, d_0$  : **náhodné** hyperparametry, tj.  $\zeta = (b_0, d_0)^\top$
- $\mu_0, k_0, b_0, p_0, r_0, d_0, g_0, h_0$  : pevné hyperparametry
- Graficky lze zpřehlednit pomocí **DAGu** (*directed acyclic graph*)
  - Kolečka: náhodné uzly
  - Čtverečky: nenáhodné uzly
  - Vyjádření (podmíněných) (ne)závislostí

# “ANOVA” hierarchický model

Apriorní rozdělení (jedna z možností)

$$p(\tau, \mu, q, b_0, d_0) = \underbrace{p(\mu | q)}_{\mathcal{N}(\mu_0, k_0^{-1} q^{-1})} \underbrace{p(q | b_0)}_{\mathcal{G}(a_0, b_0)} \underbrace{p(b_0)}_{\mathcal{G}(p_0, r_0)} \underbrace{p(\tau | d_0)}_{\mathcal{G}(c_0, d_0)} \underbrace{p(d_0)}_{\mathcal{G}(g_0, h_0)}$$

**Nenáhodné (pevné) (hyper)parametry:**

- $\mu_0, k_0$
- $a_0, p_0, r_0$
- $c_0, g_0, h_0$

▣▶ Jak je volit?

# “ANOVA” hierarchický model

## Volba nenáhodných (hyper)parametrů

- Není-li k dispozici žádná rozumná apriorní informace, je snaha volit nenáhodné (hyper)parametry tak, aby výsledné apriorní rozdělení bylo co nejméně informativní (*weakly informative*).
  - $p(\psi | \mathbf{y}) \propto L_F(\psi) p(\psi)$ ,
  - $p(\psi)$  představuje “nová” pozorování ve věrohodnosti.
- Snaha volit apriorní rozdělení tak, aby vliv těchto “nových” umělých pozorování na věrohodnost byl co možná nejmenší.
  - Snaha, aby co možná nejvíce platilo  $p(\psi) \propto 1$  (viz též první část semestru).
  - Stačí, aby platilo relativně k  $L_F(\psi)$ .
- ▣ Konkrétní volba pevných hyperparametrů je často (částečně) motivována pozorovanými daty.



# “ANOVA” hierarchický model

---

(Částečně) datově motivovaná volba pevných hyperparametrů

- $\mu_0$  : apriorní střední hodnota pro  $\mathbb{E} Y_{i,j} = \mathbb{E} b_i = \mu$ 
  - $\mu_0 \approx \bar{y}$
- $k_0$  : ovlivňuje apriorní inverzní rozptyl (přesnost) pro  $\mu$ 
  - $k_0$  blízké 0
- $a_0, c_0, p_0, g_0$  : “stupně volnosti” gama rozdělení
  - obvykle mezi 0 a 1
- $r_0, h_0$  : “rate” parametr gama rozdělení v poslední hierarchické úrovni
  - obvykle se volí blízké 0

# “ANOVA” hierarchický model

---

## Aposteriorní rozdělení

- $p(\theta | \mathbf{y})$  odvodíme standardním způsobem
  - ▣▶ Jak?
- V tomto konkrétním případě lze při volbě konjugovaného systému ještě vše odvodit analyticky.
- Pro mnohé jiné volby apriorních rozdělení již analyticky odvodit nelze (problém spočítat integrál ve jmenovateli Bayesovy věty).
- Inference pomocí aposteriorního rozdělení obvykle založená na počítačové simulaci (**Monte Carlo metody**).

# 4

**(Zobecněné) lineární smíšené  
modely**

## Oddíl 4.1

# Lineární smíšený model

- Normální lineární smíšený model

$$\mathbf{Y}_i = \mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \mathbf{b}_i + \boldsymbol{\varepsilon}_i, \quad i = 1, \dots, N$$

$$\mathbf{b}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_q(\boldsymbol{\mu}, \mathbb{D})$$

$$\boldsymbol{\varepsilon}_i \sim \mathcal{N}_{n_i}(\mathbf{0}, \boldsymbol{\Sigma}_i)$$

- $\mathbf{b}_1, \dots, \mathbf{b}_N, \boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_N$  vzájemně nezávislé
-

# Normální lineární smíšený model

## Hierarchický zápis

Normální lineární smíšený model zapsaný hierarchicky

$$\mathbf{Y}_i | \mathbf{b}_i \sim \mathcal{N}_{n_i}(\mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \mathbf{b}_i, \boldsymbol{\Sigma}_i), \quad i = 1, \dots, N$$
$$\mathbf{b}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_q(\boldsymbol{\mu}, \mathbb{D})$$

- $\mathbf{Y}_1, \dots, \mathbf{Y}_N$  vzájemně nezávislé
- $\mathbf{b}_1, \dots, \mathbf{b}_N$  vzájemně nezávislé

## Oddíl 4.2

# Příklady

## Příklad: NTP TER84073 pokus na myších

### Normální lineární smíšený model

- Praviděpodobnostní reprezentace dat:

$Y_{i,j}$  = hmotnost  $j$ -tého zárodku  $i$ -té myši,  $i = 1, \dots, N$ ,  $j = 1, \dots, n_i$

- Normální LMM:

$$Y_{i,j} | b_i \sim \mathcal{N}(b_i, \sigma^2)$$

$$\mathbf{Y}_i | b_i \sim \mathcal{N}_{n_i}(b_i \mathbf{1}, \boldsymbol{\Sigma}_i)$$

$$b_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, d^2)$$

- V obecném značení

$$\mathbb{X}_i \text{ není, } \mathbb{Z}_i = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, \quad \boldsymbol{\Sigma}_i = \sigma^2 \mathbf{I}_{n_i}, \quad \mathbb{D} = d^2$$

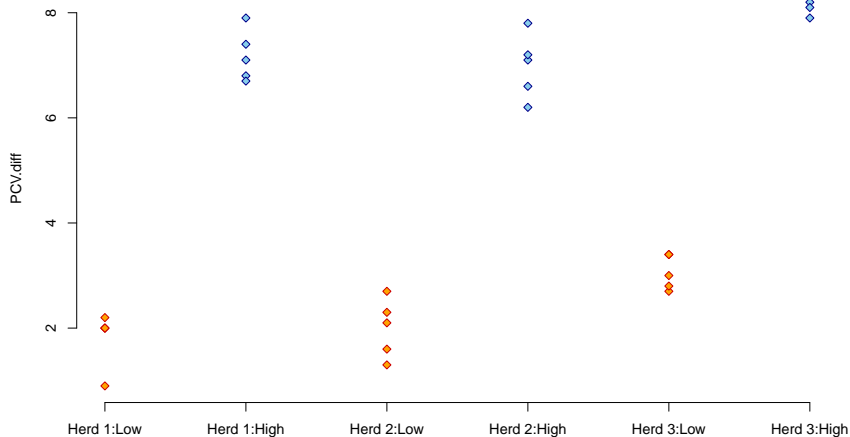


## Příklad: Vliv dávky Berenilu na léčbu trypanosomosis

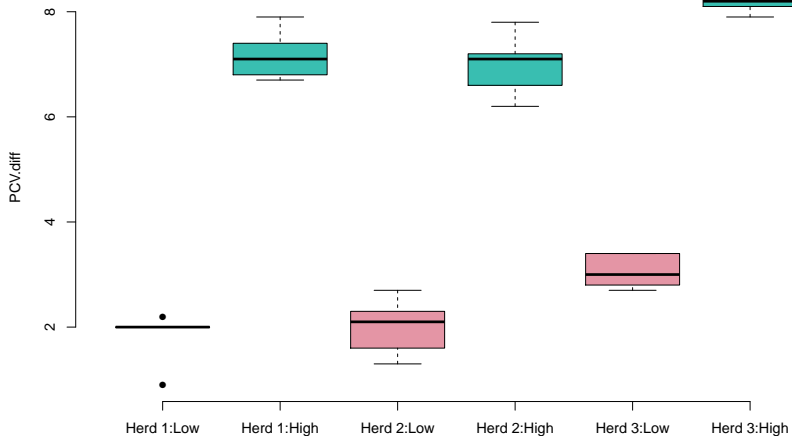
- Experiment mající za úkol vyhodnotit vliv dávky léku Berenil na úspěšnost léčby trypanosomosis u krav
- Úspěšnost léčby je měřena pomocí přírůstku PCV (packed cell volume) po určité době podávání léku
- Zajímá nás, zda existuje rozdíl v úspěšnosti léčby nízkou či vysokou dávkou Berenilu
- Experiment proběhl na kravách ze 3 stád (z různých farem)
- Každé stádo bylo rozděleno náhodně na dvě (přibližně stejné) části, přičemž krávy z jedné části stáda byly ošetřovány nízkou dávkou Berenilu, krávy z druhé části stáda byly ošetřovány vysokou dávkou Berenilu
- **Pravděpodobnostní reprezentace dat:**

$Y_{i,j}$  = PCV přírůstek u  $j$ -té krávy  $i$ -tého stáda

## Příklad: Vliv dávky Berenilu na léčbu trypanosomosis



## Příklad: Vliv dávky Berenilu na léčbu trypanosomosis



# Příklad: Vliv dávky Berenilu na léčbu trypanosomosis

## Lineární smíšený model

- Jestliže lze předpokládat, že dávka má stejný účinek ve všech stádech, lze výsledek experimentu reprezentovat následujícím LMM

$$Y_{i,j} = b_i + x_{i,j} \beta + \varepsilon_{i,j}$$

- $x_{i,j} \begin{cases} 0, & \text{jestliže } (i,j)\text{-tá kráva léčena nízkou dávkou} \\ 1, & \text{jestliže } (i,j)\text{-tá kráva léčena vysokou dávkou} \end{cases}$
- $b_i$  : střední účinek nízké dávky v  $i$ -tém stádu
  - náhodný efekt stáda
  - nepředpokládáme, že je stejné u všech stád
  - předpokládáme, že stáda jsou náhodně vybrána z populace stád
  - $\mathbb{E}b_i = \mu =$  střední účinek nízké dávky v populaci
- $\beta$  : rozdíl mezi účinky vysoké a nízké dávky
  - konstantní (pevný efekt)
  - předpokládáme, že je stejné u všech stád
  - populační rozdíl mezi účinky vysoké a nízké dávky
- $\varepsilon_{i,j}$  : náhodná odchylka  $(i,j)$ -té krávy od střední hodnoty její části (vysoká/nízká dávka) jejího stáda

# Příklad: Vliv dávky Berenilu na léčbu trypanosomosis

## Lineární smíšený model

$$Y_{i,j} = b_i + x_{i,j} \beta + \varepsilon_{i,j}$$

$$\mathbb{E}b_i = \mu$$

- **Parametr hlavního zájmu:**  $\beta$   
≡ populační rozdíl v efektu léčby vysokou a nízkou dávkou

## Příklad: Vliv dávky Berenilu na léčbu trypanosomosis

### Normální lineární smíšený model

- Normální LMM:

$$Y_{i,j} | b_i \sim \mathcal{N}(b_i + x_{i,j}\beta, \sigma^2)$$

$$\mathbf{Y}_i | b_i \sim \mathcal{N}_{n_i}(\mathbf{Z}_i b_i + \mathbb{X}_i \beta, \sigma^2 \mathbf{I}_{n_i})$$

$$b_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, d^2)$$

- V obecném značení

$$\mathbb{X}_i = \begin{pmatrix} x_{i,1} \\ \vdots \\ x_{i,n_i} \end{pmatrix}, \quad = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ \vdots \\ 1 \end{pmatrix}, \quad \mathbf{Z}_i = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, \quad \boldsymbol{\Sigma}_i = \sigma^2 \mathbf{I}_{n_i}, \quad \mathbb{D} = d^2$$

## Příklad: Vliv dávky Berenilu na léčbu trypanosomosis

### Lineární smíšený model 2

- Jestliže nelze předpokládat, že dávka má stejný účinek ve všech stádech, lze výsledek experimentu reprezentovat následujícím LMM

$$Y_{i,j} = b_{i,1} + z_{i,j} b_{i,2} + \varepsilon_{i,j}$$

- $z_{i,j} \begin{cases} 0, & \text{jestliže } (i,j)\text{-tá kráva léčena nízkou dávkou} \\ 1, & \text{jestliže } (i,j)\text{-tá kráva léčena vysokou dávkou} \end{cases}$
- $b_{i,1}$  : střední účinek nízké dávky v  $i$ -tém stádu
  - náhodný efekt stáda (nepředpokládáme, že je stejné u všech stád)
  - předpokládáme, že stáda jsou náhodně vybrána z populace stád
  - $\mathbb{E}b_{i,1} = \mu_1 =$  střední účinek nízké dávky v populaci
- $b_{i,2}$  : rozdíl mezi účinky vysoké a nízké dávky v  $i$ -tém stádu
  - náhodný efekt stáda (nepředpokládáme, že je stejné u všech stád)
  - $\mathbb{E}b_{i,2} = \mu_2 =$  střední rozdíl mezi účinky vysoké a nízké dávky v populaci
- $\varepsilon_{i,j}$  : náhodná odchylka  $(i,j)$ -té krávy od střední hodnoty její části (vysoká/nízká dávka) jejího stáda

## Příklad: Vliv dávky Berenilu na léčbu trypanosomosis

### Lineární smíšený model 2

$$Y_{i,j} = b_{i,1} + b_{i,2} z_{i,j} + \varepsilon_{i,j}$$

$$\mathbb{E}b_{i,1} = \mu_1$$

$$\mathbb{E}b_{i,2} = \mu_2$$

- **Parametr hlavního zájmu:**  $\mu_2$   
≡ populační rozdíl v efektu léčby vysokou a nízkou dávkou



## Příklad: Vliv dávky Berenilu na léčbu trypanosomosis

### Normální lineární smíšený model 2

- Označme  $\mathbf{b}_i = (b_{i,1}, b_{i,2})$
- Normální LMM:

$$Y_{i,j} | \mathbf{b}_i \sim \mathcal{N}(b_{i,1} + z_{i,j} b_{i,2}, \sigma^2)$$

$$\mathbf{Y}_i | \mathbf{b}_i \sim \mathcal{N}_{n_i}(\mathbf{Z}_i \mathbf{b}_i, \boldsymbol{\Sigma}_i)$$

$$\mathbf{b}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_2(\boldsymbol{\mu}, \mathbb{D})$$

- V obecném značení

$$\mathbb{X}_i \text{ není, } \mathbf{Z}_i = \begin{pmatrix} 1 & 0 \\ \vdots & \vdots \\ 1 & 0 \\ 1 & 1 \\ \vdots & \vdots \\ 1 & 1 \end{pmatrix}, \quad \boldsymbol{\Sigma}_i = \sigma^2 \mathbf{I}_{n_i}, \quad \boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \quad \mathbb{D} = \begin{pmatrix} d_1^2 & d_{1,2} \\ d_{1,2} & d_2^2 \end{pmatrix}$$

## Příklad: Vývoj hladiny bilirubinu u pacientů s PBC

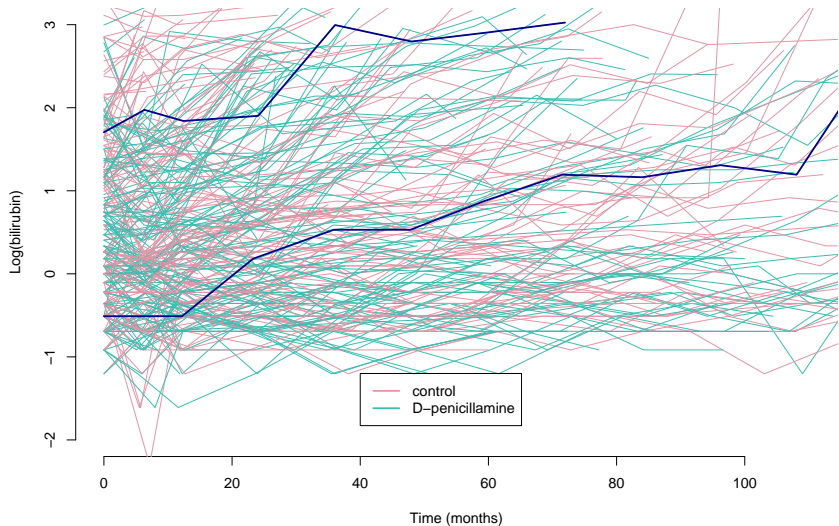
---

### Longitudinální studie

- Studie, která proběhla v letech 1974–1984 na Mayo Clinic
- 312 pacientů s PBC (primary biliary cirrhosis)
- 158 pacientů léčeno D-penicillaminem
- 154 pacientů léčeno pouze standardní léčbou
- Jedním z cílů studie je srovnat skupiny (D-penicillamin vs. kontrolní) vzhledem k vývoji hladiny bilirubinu (jeden z ukazatelů vážnosti PBC)
- Pacienti chodili v zadaných (ne zcela pravidelných) intervalech na vyšetření, kde byla zjišťována (mimo jiného) hladina bilirubinu
- Medián doby sledování byl 6,3 roku (IQR 3,7 – 8,9 roku)
- **Pravděpodobnostní reprezentace dat:**

$Y_{i,j}$  = logaritmus hladiny bilirubinu  $i$ -tého pacienta v čase  $t_{i,j}$

## Příklad: Vývoj hladiny bilirubinu u pacientů s PBC



## Příklad: Vývoj hladiny bilirubinu u pacientů s PBC

### Možný lineární smíšený model

- Vývoj hladiny logaritmického bilirubinu v čase se zdá být u každého pacienta lineární
- Možný model pro jednoho pacienta: **přímka v čase**
- Každý pacient však začíná v čase 0 na jiné úrovni
  - různé absolutní členy (intercepty) pro různé pacienty
- Růst/pokles log-bilirubinu může být jinak rychlý u jednotlivých pacientů
  - různé směrnice pro různé pacienty

III ▶  $Y_{i,j} = b_{i,1} + b_{i,2} t_{i,j} + \varepsilon_{i,j}$

- $b_{i,1}$  : absolutní člen přímky ( $\equiv$  trend)  $i$ -tého pacienta
- $b_{i,2}$  : směrnice přímky  $i$ -tého pacienta
- $\varepsilon_{i,j}$  : náhodná odchylka od celkového trendu v čase  $t_{i,j}$

## Příklad: Vývoj hladiny bilirubinu u pacientů s PBC

### Možný lineární smíšený model

$$Y_{i,j} = b_{i,1} + b_{i,2} t_{i,j} + \varepsilon_{i,j}$$

- $\mathbf{b}_i = (b_{i,1}, b_{i,2})^\top$  : náhodný efekt

▮▮▮▮▶ reprezentuje výběr z populace pacientů

- $\mathbb{E}b_{i,1} = \mu_1$  : střední (populační) absolutní člen

- $\mathbb{E}b_{i,2} = \mu_2$  : střední (populační) směrnice

▮▮▮▮▶  $\mathbb{E}Y_{i,j} = \mu_1 + \mu_2 t_{i,j}$

≡ střední vývoj log-bilirubinu v čase v populaci

## Příklad: Vývoj hladiny bilirubinu u pacientů s PBC

### Zahrnutí informace o skupině (kontrolní/D-penicillamine)

- Pacienti byli na začátku studie přiřazováni do skupin náhodně
  - ▮ Lze předpokládat stejnou střední hladinu log-bilirubinu v obou skupinách
  - ▮ Stejný střední absolutní člen v obou skupinách
- D-penicillamine může vést k jinak rychlé změně hladiny log-bilirubinu
  - ▮ Potřeba umožnit různé střední směrnice v jednotlivých skupinách

- Necht'  $x_i$   $\begin{cases} 0, & \text{jestliže } i\text{-tý pacient patří do kontrolní skupiny} \\ 1, & \text{jestliže } i\text{-tý pacient patří do D-penicillamine skupiny} \end{cases}$

- ▮ Model umožňující různé střední směrnice v jednotlivých skupinách

$$Y_{i,j} = b_{i,1} + b_{i,2} t_{i,j} + \beta t_{i,j} x_{i,j} + \varepsilon_{i,j}$$

## Příklad: Vývoj hladiny bilirubinu u pacientů s PBC

Model zahrnující informaci o skupině (kontrolní/D-penicillamine)

$$Y_{i,j} = b_{i,1} + b_{i,2} t_{i,j} + \beta t_{i,j} x_{i,j} + \varepsilon_{i,j}$$

- $\mathbf{b}_i = (b_{i,1}, b_{i,2})^\top$  : náhodný efekt

▸ reprezentuje výběr z populace pacientů

- $\mathbb{E}b_{i,1} = \mu_1$  : střední (populační) absolutní člen

▸ Při  $x_i = 0$  je  $\mathbb{E}Y_{i,j} = \mu_1 + \mu_2 t_{i,j}$

≡ střední vývoj log-bilirubinu v čase v kontrolní skupině

▸ Při  $x_i = 1$  je  $\mathbb{E}Y_{i,j} = \mu_1 + (\mu_2 + \beta) t_{i,j}$

≡ střední vývoj log-bilirubinu v čase ve skupině D-penicillamine

- $\mathbb{E}b_{i,2} = \mu_2$  : střední (populační) směrnice v kontrolní skupině

- $\mathbb{E}(b_{i,2} + \beta) = \mu_2 + \beta$  : střední (populační) směrnice ve skupině D-penicillamine

▸  $\beta$  střední (populační) rozdíl směrnic mezi skupinou D-penicillamine a kontrolní skupinou

## Příklad: Vývoj hladiny bilirubinu u pacientů s PBC

---

### Normální lineární smíšený model

- Normální LMM:

$$Y_{i,j} | \mathbf{b}_i \sim \mathcal{N}(b_{i,1} + b_{i,2} t_{i,j} + \beta t_{i,j} x_{i,j}, \sigma^2)$$

$$\mathbf{Y}_i | \mathbf{b}_i \sim \mathcal{N}_{n_i}(\mathbb{X}_i \beta + \mathbb{Z}_i \mathbf{b}_i, \boldsymbol{\Sigma}_i)$$

$$\mathbf{b}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_2(\boldsymbol{\mu}, \mathbb{D})$$

- 
- Jak vypadají  $\mathbb{X}_i$ ,  $\mathbb{Z}_i$ ,  $\boldsymbol{\Sigma}_i$ ?



## Oddíl 4.3

# Lineární smíšený model bayesovsky

## Zápis 1

$$\mathbf{Y}_i = \mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \mathbf{b}_i + \boldsymbol{\varepsilon}_i, \quad i = 1, \dots, N$$

$$\mathbf{b}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_q(\boldsymbol{\mu}, \mathbb{D})$$

$$\boldsymbol{\varepsilon}_i \sim \mathcal{N}_{n_i}(\mathbf{0}, \boldsymbol{\Sigma}_i)$$

- $\mathbf{b}_1, \dots, \mathbf{b}_N, \boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_N$  vzájemně nezávislé

## Zápis 2

$$\mathbf{Y}_i | \mathbf{b}_i \sim \mathcal{N}_{n_i}(\mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \mathbf{b}_i, \boldsymbol{\Sigma}_i), \quad i = 1, \dots, N$$

$$\mathbf{b}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_q(\boldsymbol{\mu}, \mathbb{D})$$

- $\mathbf{Y}_1, \dots, \mathbf{Y}_N$  vzájemně nezávislé
- $\mathbf{b}_1, \dots, \mathbf{b}_N$  vzájemně nezávislé

# Normální lineární smíšený model

Parametry ve frekventistickém pojetí

“Čisté” parametry (parametry ve frekventistickém smyslu)

- $\beta$  : pevné efekty
  - (populační) vliv regresorů zahrnutých v matici  $\mathbb{X}$  na odezvu
- $\mu = \mathbb{E}\mathbf{b}_i$  ( $i = 1, \dots, N$ ) : střední hodnoty náhodných efektů
  - (populační) vliv regresorů zahrnutých v matici  $\mathbb{Z}$  na odezvu
- $\Sigma_i = \text{var}(\mathbf{Y}_i | \mathbf{b}_i)$  ( $i = 1, \dots, N$ ) : “vnitroskupinová” varianční matice
  - často se předpokládá  $\Sigma_i = \sigma^2 \mathbf{I}_{n_i}$ 
    - ▮▮▮▮ podmíněná nezávislost
  - pomocí jiné struktury pro matici  $\Sigma_i$  lze modelovat i jiné struktury závislosti (AR(d), ...)
- $\mathbb{D} = \text{var}\mathbf{b}_i$  ( $i = 1, \dots, N$ ) : “meziskupinová” varianční matice
  - obvykle se nepředpokládá žádná speciální struktura  $\mathbb{D}$  a pouze se požaduje, aby  $\mathbb{D} > 0$  (pozitivně definitní matice)


# Normální lineární smíšený model

## Frekventistická věrohodnost

- **Frekventistické parametry:**  $\psi = (\beta^\top, \mu^\top, \underbrace{\text{par}(\Sigma_1, \dots, \Sigma_N), \text{par}(\mathbb{D})}^\top)^\top$   
často pouze  $\sigma^2$
- **Frekventistická věrohodnost:**

$$L_F(\psi) = p(\mathbf{y} | \psi) = \dots = \prod_{i=1}^N \varphi(\mathbf{y}_i | \mathbb{X}_i \beta + \mathbb{Z}_i \mu, \mathbb{V}_i),$$

$$\text{kde } \mathbb{V}_i = \mathbb{Z}_i \mathbb{D} \mathbb{Z}_i^\top + \Sigma_i$$

- Při odhadu metodou **maximální věrohodnosti** potřeba maximalizovat  $L_F(\psi)$  při omezeních  $\Sigma_i > 0$  (pro všechna  $i = 1, \dots, N$ ) a  $\mathbb{D} > 0$ 
  -  balíčky `lme4`, `nlme`
  - SAS procedura `MIXED`

# Normální lineární smíšený model

---

## Skrytá data a věrohodnost bayesovského modelu

- Další náhodné složky modelu:
  - ≡ vektory náhodných efektů  $\mathbf{b}_1, \dots, \mathbf{b}_N$
  - ▮▶ Další “parametry” při bayesovském přístupu
  - ≡ skrytá data
- **Skrytá data:**  $\xi = (\mathbf{b}_1^\top, \dots, \mathbf{b}_N^\top)^\top$
- **Věrohodnost bayesovského modelu:**

$$L(\xi, \psi) = p(\mathbf{y} | \xi, \psi) = \dots = \prod_{i=1}^N \varphi(\mathbf{y}_i | \mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \mathbf{b}_i, \boldsymbol{\Sigma}_i)$$

## Oddíl 4.4

# Zobecněný lineární smíšený model

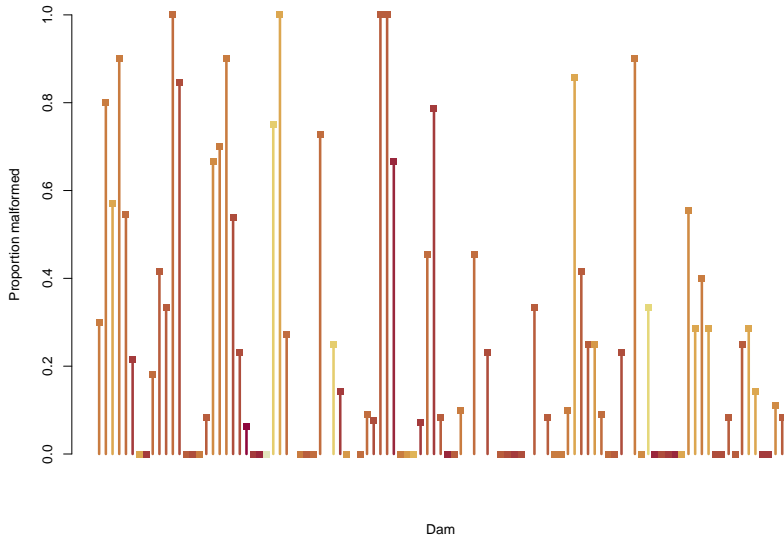
- Data z *National Toxicology Program*.
- 94 těhotným myším byla podána v předem určených momentech určená množství etylenglykolu (EG).
- V 17. dnu těhotenství byly myši usmrceny a následně byl zaznamenán počet zárodků ( $n_i, i = 1, \dots, 94$ ) a indikace, zda se u jednotlivých zárodků vyskytovala vývojová vada.
- **Pravděpodobnostní reprezentace dat:**

$$Y_{i,j} = \begin{cases} 0, & \text{jestliže } j\text{-tý zárodek } i\text{-té myši bez vývojové vady,} \\ 1, & \text{jestliže } j\text{-tý zárodek } i\text{-té myši s vývojovou vadou.} \end{cases}$$

- **Primární cíl:**  
Odhad a inference pro  $\pi = \mathbb{E}Y_{i,j} = P(Y_{i,j} = 1)$ .

# Příklad: NTP TER84073 pokus na myších

Pozorované proporce zárodků s vývojovými vadami





# Příklad: NTP TER84073 pokus na myších

## Možný model

**Možný model** ( $j = 1, \dots, n_i$  pro každé  $i = 1, \dots, N$ )

- $Y_{i,j} | \pi_i \sim \mathcal{A}(\pi_i)$ .

- $\pi_i \stackrel{\text{i.i.d.}}{\sim}$  z nějakého rozdělení.

► Reprezentuje fakt, že každá myš má jiné (genetické apod.) dispozice k tomu, aby se u jejích zárodků vyvinuly vývojové vady.

- Myši zahrnuté do studie jsou náhodným výběrem z populace myší a proto je rozumné předpokládat, že též  $\pi_i$  ( $i = 1, \dots, N$ ) jsou náhodné.

- Abychom se nemuseli starat o omezení  $0 < \pi_i < 1$  ( $i = 1, \dots, N$ ), bude užitečné použít vhodnou reparametrizaci, např.

$$\pi_i = \frac{e^{b_i}}{1 + e^{b_i}}, \quad b_i = \text{logit}(\pi_i) = \log\left(\frac{\pi_i}{1 - \pi_i}\right)$$

a v modelu uvažovat  $b_i \stackrel{\text{i.i.d.}}{\sim}$  z nějakého rozdělení.

► **Logistická regrese s náhodnými efekty.**

# Příklad: NTP TER84073 pokus na myších

Logistická regrese s normálně rozdělenými náhodnými efekty

**Možný model** ( $j = 1, \dots, n_i$  pro každé  $i = 1, \dots, N$ )

- $Y_{i,j} | b_i$  nezávislé s rozdělením  $\mathcal{A}(\pi_i)$ , kde  $\pi_i = \frac{e^{b_i}}{1+e^{b_i}}$ .
- $b_i$  i.i.d. s rozdělením  $\mathcal{N}(\mu, d^2)$ .

**Parametry** (klasické)

- $\psi = (\mu, d^2)^\top$ .

**Věrohodnost** (frekventistická)

$$\begin{aligned} L_F(\psi) &= p(\mathbf{y} | \psi) = \prod_{i=1}^N p(\mathbf{y}_i | \psi) = \prod_{i=1}^N \int_{\mathbb{R}} \prod_{j=1}^{n_i} p(y_{i,j} | b_i, \psi) p(b_i | \psi) db_i \\ &= \prod_{i=1}^N \int_{\mathbb{R}} \pi_i^{\sum_{j=1}^{n_i} y_{i,j}} (1 - \pi_i)^{n_i - \sum_{j=1}^{n_i} y_{i,j}} \varphi(b_i | \mu, d^2) db_i \\ &= \prod_{i=1}^N \int_{\mathbb{R}} \frac{e^{b_i \sum_{j=1}^{n_i} y_{i,j}}}{(1 + e^{b_i})^{n_i}} \varphi(b_i | \mu, d^2) db_i. \end{aligned}$$

## Oddíl 4.5

# Apriorní rozdělení

## Apriorní rozdělení

- Využijeme rozklad, který sleduje hierarchickou strukturu modelu

$$p(\xi, \psi) = p(\xi | \psi) p(\psi)$$

- První část

$$p(\xi | \psi) = p(\mathbf{b}_1, \dots, \mathbf{b}_N | \beta, \mu, \Sigma_1, \dots, \Sigma_N, \mathbb{D}) =$$

$$p(\mathbf{b}_1, \dots, \mathbf{b}_N | \mu, \mathbb{D}) = \dots = \prod_{i=1}^N \varphi(\mathbf{b}_i | \mu, \mathbb{D})$$

- Druhá část  $p(\psi)$

- “standardní” apriorní rozdělení pro “čisté” parametry
- obvykle se při jeho specifikaci zavádějí na principu obecných hierarchických modelů další náhodné hyperparametry  $\zeta$  (snaha zabránit nadměrnému vlivu zvoleného apriorního rozdělení na rozdělení aposteriorní)

# Normální (zobecněný) lineární smíšený model

## Apriorní rozdělení

- **Parametry bayesovského modelu**  $\theta = (\xi, \psi^\top, \zeta^\top)^\top$

- $\xi = (\mathbf{b}_1^\top, \dots, \mathbf{b}_N^\top)^\top$

- $\psi = (\beta^\top, \mu^\top, \text{par}(\Sigma_1, \dots, \Sigma_N), \text{par}(\mathbb{D}))^\top$

- $\zeta$  : případné další hyperparametry

- **Rozklad apriorního rozdělení:**

$$p(\theta) = p(\xi, \psi, \zeta) = p(\xi | \psi) p(\psi | \zeta) p(\zeta)$$

$$= \left\{ \prod_{i=1}^N \varphi(\mathbf{b}_i | \mu, \mathbb{D}) \right\} p(\psi | \zeta) p(\zeta)$$

# Normální (zobecněný) lineární smíšený model

## Možné volby pro nestrukturální část apriorního rozdělení

- Potřeba zvolit  $p(\psi)$ , které často specifikujeme hierarchicky pomocí dalších hyperparametrů jako

$$p(\psi) = \int p(\psi | \zeta) p(\zeta) d\zeta,$$

to jest  $p(\psi, \zeta) = p(\psi | \zeta) p(\zeta)$

- $\psi = (\beta^\top, \mu^\top, \text{par}(\Sigma_1, \dots, \Sigma_N), \text{par}(\mathbb{D}))^\top$
- V dalším se budeme zabývat pouze situací, kdy  $\Sigma_i = \sigma^2 \mathbf{I}_{n_i}$ , tj.  $\text{par}(\Sigma_1, \dots, \Sigma_N) = \sigma^2$
- Označíme dále

$$\tau = \sigma^{-2}$$

$$\mathbb{Q} = \mathbb{D}^{-1}$$

$$\Rightarrow \psi = (\beta^\top, \mu^\top, \tau, \text{par}(\mathbb{Q}))^\top$$

# Normální (zobecněný) lineární smíšený model

## Možné volby pro nestrukturální část apriorního rozdělení

- Obvykle se předpokládá apriorní nezávislost pro jednotlivé sady “příbuzných” parametrů, např.

$$p(\beta, \mu, \tau, \mathbb{Q}) = p(\beta) p(\mu) p(\tau) p(\mathbb{Q}),$$

respektive

$$p(\beta, \mu, \tau, \mathbb{Q} | \zeta) = p(\beta | \zeta^{(1)}) p(\mu | \zeta^{(2)}) p(\tau | \zeta^{(3)}) p(\mathbb{Q} | \zeta^{(4)}),$$

---

kde  $\zeta = (\zeta^{(1)\top}, \zeta^{(2)\top}, \zeta^{(3)\top}, \zeta^{(4)\top})^\top$

# Normální (zobecněný) lineární smíšený model

Možné volby pro nestrukturální část apriorního rozdělení

$\beta$ ,  $\mu$  mají interpretaci středních hodnot

- Smysluplné apriorní rozdělení:

$$p(\beta) \propto 1$$

$$p(\mu) \propto 1$$

- Jiné smysluplné apriorní rozdělení:

$$p(\beta) \sim \mathcal{N}(\beta_0, \Sigma_0^\beta)$$

$$p(\mu) \sim \mathcal{N}(\mu_0, \Sigma_0^\mu),$$

- $\beta_0$ ,  $\Sigma_0^\beta$ ,  $\mu_0$ ,  $\Sigma_0^\mu$  : pevné/náhodné hyperparametry
- častá smysluplná volba:
  - $\beta_0 = \mathbf{0}$  (kromě abs. členu modelu)
  - $\mu_0 = \mathbf{0}$  (kromě abs. členu modelu)
  - $\Sigma_0^\beta, \Sigma_0^\mu$  : diagonální matice s velkými (co je velké?) čísly na diagonále



# Normální (zobecněný) lineární smíšený model

## Možné volby pro nestrukturální část apriorního rozdělení

$\tau$  je inverzní rozptyl odchylek jednotlivých pozorování  $i$ -tého subjektu od střední úrovně (závislé na regresorech) tohoto subjektu

- Smysluplné apriorní rozdělení:

$$p(\tau) \sim \mathcal{G}(c_\tau, d_\tau)$$

- $c_\tau, d_\tau$  : pevné/náhodné hyperparametry
- mimo jiné zajišťuje, že  $P(\tau > 0 | \mathbf{Y}) = 1$
- častá smysluplná volba hyperparametrů:
  - $c_\tau$  : apriorní “stupně volnosti”, tj.  $c_\tau \in (0, 1]$  vede k slabě informativnímu rozdělení
  - pro připomenutí:  $\mathbb{E}_\tau = \frac{c_\tau}{d_\tau}$ ,  $\text{var}_\tau = \frac{c_\tau}{d_\tau^2}$ 
    - ▮  $d_\tau$  : “přesnost” apriorního gama rozdělení
  - $d_\tau$  blízké 0 může vést k slabě informativnímu rozdělení
  - volba  $d_\tau$  však může poměrně značně ovlivnit tvar aposteriorního rozdělení
    - ▮  $d_\tau$  se často volí jako **náhodné** s dalším (již pevně zvoleným) gama rozdělením jako apriorním

# Normální (zobecněný) lineární smíšený model

---

Možné volby pro nestrukturální část apriorního rozdělení

Q je inverzní varianční matice “úrovni” jednotlivých subjektů

- Potřeba **vícerozměrné** apriorní rozdělení, které s pravděpodobností 1 vygeneruje pozitivně definitní matici

▣ **Wishartovo** rozdělení

$\mathbb{Q} \sim \mathcal{W}_p(\nu, \Xi)$ , jestliže

$$\mathbb{Q} = \sum_{i=1}^{\nu} \mathbf{z}_i \mathbf{z}_i^{\top},$$

- kde
- $\mathbf{z}_1, \dots, \mathbf{z}_\nu \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_p(\mathbf{0}, \Xi)$
  - $\Xi$  je pozitivně definitní měřítková matice
  - $\nu > p - 1$  jsou stupně volnosti
  - Zřejmě platí:  $P(\mathbb{Q} > 0) = 1$
  - Jedná se o vícerozměrné rozšíření  $\chi_\nu^2$  rozdělení

$\mathbb{Q} \sim \mathcal{W}_p(\nu, \Xi)$  má hustotu

$$p(\mathbb{Q}) = \left\{ 2^{\frac{\nu p}{2}} \pi^{\frac{p(p-1)}{4}} \prod_{i=1}^p \Gamma\left(\frac{\nu+1-i}{2}\right) \right\}^{-1} |\Xi|^{-\frac{\nu}{2}}$$

$$|\mathbb{Q}|^{\frac{\nu-p-1}{2}} \exp\left\{-\frac{1}{2} \text{tr}(\Xi^{-1} \mathbb{Q})\right\}, \quad \mathbb{Q} \text{ pozitivně definitní}$$

- $\nu > p - 1$  : “stupně volnosti” lze uvažovat i neceločíselné
  - ▮ zobecnění klasického Wishartova rozdělení
- $\mathbb{E}\mathbb{Q} = \nu\Xi$
- $\mathcal{W}_1(\nu, 1) \equiv \mathcal{G}\left(\frac{\nu}{2}, \frac{1}{2}\right) \equiv \chi_\nu^2$
- $\mathcal{W}_1(\nu, \Xi) \equiv \mathcal{G}\left(\frac{\nu}{2}, \frac{\Xi^{-1}}{2}\right)$

# Normální (zobecněný) lineární smíšený model

Možné volby pro nestrukturální část apriorního rozdělení

$\mathbb{Q}$  je inverzní varianční matice “úrovní” jednotlivých subjektů

- Smysluplné apriorní rozdělení:

$$p(\mathbb{Q}) \sim \mathcal{W}(\nu_Q, \Xi_Q)$$

- $\Xi_Q, \nu_Q$  : pevné/náhodné hyperparametry
- častá smysluplná volba hyperparametrů:
  - $\nu_Q$  : apriorní “stupně volnosti”, tj.  $\nu_\tau \in (p - 1, p]$  vede k slabě informativnímu rozdělení
  - $\Xi_Q$  se obvykle volí jako **diagonální** matice, např.  $\Xi_Q = \text{diag}(\gamma_{Q,1}^{-1}, \dots, \gamma_{Q,p}^{-1})$
  - inverze měřítkové matice ( $\Xi_Q^{-1}$ ) je “přesností” Wishartova rozdělení
    - ▮▮▮ diagonální hodnoty  $\Xi_Q^{-1}$  (tj.  $\gamma_{Q,1}, \dots, \gamma_{Q,p}$ ) blízké 0 mohou vést k slabě informativnímu apriornímu rozdělení
  - volby  $\gamma_{Q,1}, \dots, \gamma_{Q,p}$  však mohou poměrně značně ovlivnit tvar aposteriorního rozdělení
    - ▮▮▮  $\gamma_{Q,1}, \dots, \gamma_{Q,p}$  se proto často volí jako **náhodné**, apriorně vzájemně nezávislé, s dalšími (již pevně zvolenými) gama rozděleními jako apriorními

# Normální (zobecněný) lineární smíšený model

## Shrnutí

- “čisté” parametry:  $\psi = (\beta^\top, \mu^\top, \tau, \text{par}(\mathbb{Q}))^\top$
- Skrytá data:  $\xi = (\mathbf{b}_1^\top, \dots, \mathbf{b}_N^\top)^\top$
- Věrohodnost bayesovského modelu:

$$L(\xi, \psi) = p(\mathbf{y} | \xi, \psi) = \dots = \prod_{i=1}^N \varphi(\mathbf{y}_i | \mathbb{X}_i \beta + \mathbb{Z}_i \mathbf{b}_i, \tau^{-1} \mathbf{I}_{n_i})$$

- Dekompozice **apriorního rozdělení**:

$$p(\xi, \psi) = p(\xi | \psi) p(\psi) = \prod_{i=1}^N \varphi(\mathbf{b}_i | \mu, \mathbb{Q}^{-1}) p(\beta) p(\mu) p(\tau) p(\mathbb{Q})$$

## Oddíl 4.6

# Aposterorní rozdělení

# Normální (zobecněný) lineární smíšený model

---

## Aposteriorní rozdělení

- Nejsou-li žádné náhodné hyperparametry:

$$p(\xi, \psi | \mathbf{y}) \propto L(\xi, \psi) p(\xi | \psi) p(\psi)$$

---

- Jsou-li některé hyperparametry (označené jako  $\zeta$ ) náhodné:

$$p(\xi, \psi, \zeta | \mathbf{y}) \propto L(\xi, \psi) p(\xi | \psi) p(\psi | \zeta) p(\zeta)$$

$$p(\xi, \psi | \mathbf{y}) \propto \int L(\xi, \psi) p(\xi | \psi) p(\psi | \zeta) p(\zeta) d\zeta$$

---



# Normální (zobecněný) lineární smíšený model

## Aposteriorní rozdělení "čistých" parametrů

- Marginální aposteriorní rozdělení "čistých" parametrů:

$$p(\psi | \mathbf{y}) \propto \int L(\xi, \psi) p(\xi | \psi) p(\psi | \zeta) p(\zeta) d(\xi, \zeta)$$

- Díky hierarchické struktuře též platí

$$p(\psi | \mathbf{y}) \propto L_F(\psi) p(\psi),$$

$$\text{neboť } p(\psi) = \int p(\psi | \zeta) p(\zeta) d\zeta$$

$$L_F(\psi) = \int L(\xi, \psi) p(\xi | \psi) d\xi$$

a platí Fubiniova věta

## Oddíl 4.7

# Inference založená na aposteriorním rozdělení

- **Sdružené** aposteriorní rozdělení pro  $\theta = (\psi^\top, \xi^\top, \zeta^\top)^\top$  vyjádříme snadno až na multiplikační konstantu
- Primárně nás však zajímají hlavně **marginální** aposteriorní rozdělení pro sady parametrů nebo dokonce jednotlivé parametry
  - $p(\beta | \mathbf{y}), p(\beta_j | \mathbf{y}), p(\mu | \mathbf{y}), p(\tau | \mathbf{y}), \dots$
  - z nich odvozujeme aposteriorní střední hodnotu, věrohodnostní množiny atp.
- Při výpočtu **marginálních** aposteriorních rozdělení se již nelze vyhnout **integrování**
  - analyticky proveditelné pro jednodušší modely s apriorními rozděleními vykazujícími alespoň nějaký stupeň konjugovanosti
  - analyticky pracné pro složitější modely i s apriorními rozděleními vykazujícími alespoň nějaký stupeň konjugovanosti
  - analyticky často neproveditelné

- Často nás zajímá též aposteriorní rozdělení pro nějakou měřitelnou funkci  $t$  – **transformaci** původních parametrů
- Příklad:
  - $Y_{i,j} | b_i \sim \mathcal{N}(b_i, \sigma^2)$ ,  $b_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, d^2)$ ,  $i = 1, \dots, N$ ,  $j = 1, \dots, n_i$
  - Bylo:  $\text{cor}(Y_{i,j}, Y_{i,k}) = d^2 / (\sigma^2 + d^2) =: \rho$
  - $\Rightarrow$  Z  $p(\mu, \sigma^2, d^2 | \mathbf{y})$  potřeba pomocí **věty o transformaci** a **integrováním** odvodit  $p(\rho | \mathbf{y})$
  - Vzpomeňte si na svoje úspěchy na tomto poli (obvykle v poměrně jednoduchých situacích) ze 3. ročníku...
- Odvození  $p(t(\theta) | \mathbf{y})$  z  $p(\theta | \mathbf{y})$  (které již samo o sobě obvykle známe až na multiplikační konstantu) je často analyticky neproveditelné

- Vše výše řečené je důvodem toho, že bayesovská statistika byla až cca do začátku 90. let 20. století relativně málo prakticky využívána, a pokud ano, tak jenom v souvislosti s poměrně jednoduchými modely
- Řešení problému s nemožností odvozovat potřebné výrazy analyticky:
  - ▣▶ **aposteriorní inference založená na simulaci**
    - k jejímu efektivnímu použití potřeba přiměřeně **výkonné** počítačlo
    - do cca začátku 90. let 20. století jenom omezeně dostupné/nedostupné (nejenom v ČSSR)

# 5

**MCMC (Markov chain Monte Carlo)  
metody**

# Oddíl **5.1**

## **Úvod**

- Ne vždy (s reálnými aplikacemi jenom velice zřídka) jsme schopni získat náhodný výběr z požadovaného (aposteriorního) rozdělení.
- Alternativa: Konstrukce **markovského řetězce**, jehož **stacionární** rozdělení = požadované rozdělení.
- Za jistých předpokladů lze vygenerovaný **dostatečně dlouhý** markovský řetězec považovat **přibližně** za **náhodný výběr** z požadovaného rozdělení a s jako takovým s ním dále pracovat při přibližném výpočtu integrálů a jiných charakteristik tohoto rozdělení.

## ▣▶ **Markov chain Monte Carlo (MCMC)**



# Markov chain Monte Carlo

---

## Připomenutí, čím se zabýváme

- Potřebujeme generovat v ideálním případě náhodný výběr z rozdělení s hustotou  $f(\theta)$  vzhledem k  $\sigma$ -konečné míře  $\lambda$ .
  - V rámci bayesovských metod budeme obvykle používat s  $f(\theta) = p(\theta | \mathbf{y})$ .
- V dalším se předpokládají znalosti z předmětu **Náhodné procesy 1 (NMSA334)**.
- Zde jste se zabývali zejména markovskými řetězci se **spočetnou** množinou stavů.

- My potřebujeme uvažovat množinu stavů, která je obecnou podmnožinou  $\mathbb{R}^k$ .
  - Mnoho výsledků pro markovské řetězce se spočtenou množinou stavů lze zobecnit na markovské řetězce s obecnou množinou stavů.
  - V rámci této přednášky nebudeme téměř nic dokazovat, neboť potřebujeme MCMC pouze jako nástroj pro aposteriorní inferenci.
  - Důkazy a mnohé další užitečné poznatky a souvislosti lze nalézt v přednášce **Metody Markov Chain Monte Carlo (NMTP539)**.
  - Existuje též nespočet knih věnujících se MCMC.
    - Robert, C. P. and Casella, G. (2004). *Monte Carlo Statistical Methods, Second Edition*. New York: Springer-Verlag.
    - Gamerman, D. and Lopes, H. F. (2006). *Markov Chain Monte Carlo: Stochastic Simulation for Bayesian Inference, Second Edition*. Boca Raton: Chapman & Hall/CRC.

- Budeme předpokládat, že  $f(\boldsymbol{\theta})$  má hustotu vzhledem k  $\sigma$ -konečné míře  $\lambda$  na měřitelném prostoru  $(\Theta, \mathcal{T})$ .
- Typicky  $\Theta \subset \mathbb{R}^k$ .
- Obdobně jako dříve budeme zkráceně psát

$$f(\boldsymbol{\theta})d\lambda(\boldsymbol{\theta}) = f(d\boldsymbol{\theta}).$$

---

## Oddíl **5.2**

# **Markovské procesy s obecnou množinou stavů**

### Definice 5.1 Markovské jádro (*Markov kernel*).

Měřitelné zobrazení  $P : \Theta \times \mathcal{T} \rightarrow [0, 1]$  se nazývá **markovské jádro** (*Markov kernel*) na  $(\Theta, \mathcal{T})$ , jestliže

1. pro každé  $T \in \mathcal{T}$  je  $P(\cdot, T)$  nezáporná měřitelná funkce na  $\Theta$ ,
2. pro každé  $\theta \in \Theta$  je  $P(\theta, \cdot)$  pravděpodobnostní míra na  $\mathcal{T}$ .

- Jedná se o zobecnění **pravděpodobností přechodu**.
- Markovské jádro určuje pravděpodobnost přechodu ze stavu  $\theta$  do stavu  $v$  v množině  $T$ .

# Markovské procesy s obecnou množinou stavů

## Přechodová hustota

- Pro každé  $\theta \in \Theta$  existuje **hustota** rozdělení  $P(\theta, \cdot)$  (vzhledem k  $\sigma$ -konečné míře  $\lambda$ )  $\equiv$  **přechodová hustota**, kterou budeme značit jako  $p(\theta, \psi)$  a zkráceně psát

$$p(\theta, \psi)d\lambda(\psi) = p(\theta, d\psi),$$

resp. 
$$P(\theta, T) = \int_T p(\theta, \psi)d\lambda(\psi) = \int_T p(\theta, d\psi).$$

- Přechodovou hustotu lze chápat též jako podmíněnou hustotu  $\psi$  za podmínky  $\theta$  (vzhledem k  $\sigma$ -konečné míře  $\lambda$ ).

# Markovské procesy s obecnou množinou stavů

## Homogenní markovský řetězec

**Definice 5.2** Homogenní markovský řetězec (*Homogeneous Markov chain*).

Řekneme, že náhodný proces  $\{\theta^{(m)} : m = 0, \dots\}$  je **homogenní markovský řetězec** (*homogeneous Markov chain*) s přechodovým jádrem (*transition kernel*)  $P$  a počátečním rozdělením  $f_0(d\theta)$ , pokud jeho konečně rozměrná rozdělení splňují pro každé  $m \in \mathbb{N}_0$  a pro všechna  $T_0, \dots, T_m \in \mathcal{T}$  podmínku

$$P(\theta^{(0)} \in T_0, \dots, \theta^{(m)} \in T_m) = \int_{T_0} \dots \int_{T_{m-1}} P(\theta^{(m-1)}, T_m) p(\theta^{(m-2)}, d\theta^{(m-1)}) \dots p(\theta^{(0)}, d\theta^{(1)}) f_0(d\theta^{(0)}).$$

- Jak chápat tuto definici?

### Věta 5.1

Nechť  $\{\boldsymbol{\theta}^{(m)} : m = 0, \dots\}$  je homogenní markovský řetězec generovaný přechodovým jádrem  $P$  a  $h$  je omezená měřitelná funkce na  $\Theta$ . Potom pro každé  $m \in \mathbb{N}_0$  platí

$$\mathbb{E} \left[ h(\boldsymbol{\theta}^{(m+1)}) \mid \boldsymbol{\theta}^{(m)}, \dots, \boldsymbol{\theta}^{(0)} \right] = \mathbb{E} \left[ h(\boldsymbol{\theta}^{(m+1)}) \mid \boldsymbol{\theta}^{(m)} \right].$$

- Pro libovolnou  $T \in \mathcal{T}$  tedy volbou  $h = \mathbb{I}_T$  platí **markovská vlastnost**:

$$P(\boldsymbol{\theta}^{(m+1)} \in T \mid \boldsymbol{\theta}^{(m)}, \dots, \boldsymbol{\theta}^{(0)}) = P(\boldsymbol{\theta}^{(m+1)} \in T \mid \boldsymbol{\theta}^{(m)}).$$



# Markovské procesy s obecnou množinou stavů

## Rozdělení markovského řetězce v čase $m + 1$

- Necht'  $\pi$  je hustota vzhledem k  $\sigma$ -konečné míře  $\lambda$  na  $(\Theta, \mathcal{T})$ , tj.  $\pi(d\theta)$  je pravděpodobnostní rozdělení.
- Pro homogenní markovský řetězec s přechodovým jádrem  $P$  zaved'me následující značení:

$$\pi P(T) = \int_{\Theta} P(\theta, T) \pi(d\theta)$$

pro libovolnou  $T \in \mathcal{T}$ .

- Při použití různého značení máme:

$$\pi P(T) = \int_{\Theta} \int_{\mathcal{T}} p(\theta, d\psi) \pi(d\theta) = \int_{\Theta} \int_{\mathcal{T}} p(\theta, \psi) \pi(\theta) d\lambda(\psi) d\lambda(\theta).$$

# Markovské procesy s obecnou množinou stavů

## Rozdělení markovského řetězce v čase $m + 1$

- Použitím Fubiniovy věty dostaneme

$$\begin{aligned}\pi P(T) &= \int_{\Theta} \int_{\mathcal{T}} p(\theta, \psi) \pi(\theta) d\lambda(\psi) d\lambda(\theta) \\ &= \int_{\mathcal{T}} \int_{\Theta} p(\theta, \psi) \pi(\theta) d\lambda(\theta) d\lambda(\psi).\end{aligned}$$

- $\pi P(d\psi)$  je opět pravděpodobnostní rozdělení na  $(\Theta, \mathcal{T})$  a má hustotu

$$\int_{\Theta} p(\theta, \psi) \pi(\theta) d\lambda(\theta) = \int_{\Theta} p(\theta, \psi) \pi(d\theta)$$

vzhledem k míře  $\lambda$ .

- Je-li  $\pi(d\theta)$  rozdělením markovského řetězce v čase  $m - 1$ , potom  $\pi P(d\psi)$  je rozdělením markovského řetězce v čase  $m$ .

### **Definice 5.3** Stacionární rozdělení (*Stationary distribution*).

Pravděpodobnostní rozdělení  $\pi(d\theta)$  se nazývá **stacionárním rozdělením** (*stationary distribution*) homogenního markovského řetězce s přechodovým jádrem  $P$ , jestliže

$$\pi P(d\theta) = \pi(d\theta),$$

to jest, jestliže pro libovolnou  $T \in \mathcal{T}$  platí

$$\begin{aligned}\pi P(T) &= \pi(T), \\ \int_{\Theta} P(\theta, T) \pi(d\theta) &= \int_T \pi(d\theta).\end{aligned}$$

### **Definice 5.4** Reversibilita (*reversibility*) markovského řetězce.

Homogenní markovský řetězec s přechodovým jádrem  $P$  se nazývá **reversibilní** (*reversible*) vzhledem k pravděpodobnostnímu rozdělení  $\pi$ , jestliže pro libovolné  $T, S \in \mathcal{T}$  platí

$$\int_T P(\theta, S)\pi(d\theta) = \int_S P(\psi, T)\pi(d\psi).$$

## Reversibilita

- Na  $\int_T P(\theta, S)\pi(d\theta)$  lze pohlížet jako na **sdužené** pravděpodobnostní rozdělení na  $(\Theta \times \Theta, \mathcal{T} \otimes \mathcal{T})$ , které množinám  $T, S \in \mathcal{T}$  přiřadí pravděpodobnost

$$Q_1(T, S) = \int_T P(\theta, S)\pi(d\theta) = \int_T \int_S p(\theta, d\psi)\pi(d\theta).$$

- Rozdělení  $Q_1$  má hustotu

$$q_1(\theta, \psi) = p(\theta, \psi)\pi(\theta)$$

vzhledem k součinové míře  $\lambda \otimes \lambda$ , zkráceně

$$q_1(d\theta, d\psi) = p(\theta, d\psi)\pi(d\theta).$$

## Reversibilita

- Na  $\int_S P(\psi, T)\pi(d\psi)$  lze pohlížet jako na **sdužené** pravděpodobnostní rozdělení na  $(\Theta \times \Theta, \mathcal{T} \otimes \mathcal{T})$ , které množinám  $T, S \in \mathcal{T}$  přiřadí pravděpodobnost

$$Q_2(S, T) = \int_S P(\psi, T)\pi(d\psi) = \int_S \int_T p(\psi, d\theta)\pi(d\psi).$$

- Rozdělení  $Q_2$  má hustotu

$$q_2(\psi, \theta) = p(\psi, \theta)\pi(\psi)$$

vzhledem k součinové míře  $\lambda \otimes \lambda$ , zkráceně

$$q_2(d\psi, d\theta) = p(\psi, d\theta)\pi(d\psi).$$

# Markovské procesy s obecnou množinou stavů

## Reversibilita

- Reversibilita vzhledem k  $\pi$  tedy znamená, že pro libovolné  $T, S \in \mathcal{T}$  platí

$$Q_1(T, S) = Q_2(S, T).$$

- Zkráceně zapsáno

$$p(\theta, d\psi)\pi(d\theta) = p(\psi, d\theta)\pi(d\psi).$$

- Nutno chápat jako rovnost rozdělení (pravděpodobnostních měř)!
- Sdružené rozdělení stavů v časech  $m$  a  $m + 1$  je stejné jako sdružené rozdělení stavů v časech  $m + 1$  a  $m$  pro libovolné  $m = 0, 1, \dots$

# Markovské procesy s obecnou množinou stavů

---

## Reversibilita a detailní podmínka rovnováhy

- $p(\theta, d\psi)$  lze chápat též jako **podmíněné** rozdělení  $\psi$  za podmínky  $\theta$  mající hustotu  $p(\theta, \psi)$  (přechodová hustota) vzhledem k  $\sigma$ -konečné míře  $\lambda$  (v argumentu  $\psi$ ).
- Podmínku reversibility lze potom pomocí hustot přepsat jako

$$p(\theta, \psi) \pi(\theta) = p(\psi, \theta) \pi(\psi) \quad \text{pro } \lambda\text{-s.v. } \theta, \psi \in \Theta.$$

---

▣ Detailní podmínka rovnováhy (*detailed balance condition*).



### Věta 5.2 .

*Je-li homogenní markovský řetězec reversibilní vzhledem k  $\pi$ , potom  $\pi$  je jeho stacionární rozdělení.*

*Důkaz.* Stačí v definici reversibility položit  $S = \Theta$ .



- To jest, reversibilita (splnění detailní podmínky rovnováhy) je **postačující** podmínkou pro stacionaritu.

## Markovské procesy s obecnou množinou stavů

Přechodové jádro  $m$ -tého řádu

**Definice 5.5** Přechodové jádro  $m$ -tého řádu ( *$m$ -step transition probability kernel*).

Uvažujme homogenní markovský řetězec s přechodovým jádrem  $P$  a položme  $P^0(\theta, T) = \delta_\theta(T)$ . Přechodové jádro  $m$ -tého řádu ( *$m$ -step transition probability kernel*) je dáno induktivně vztahem

$$P^m(\theta, T) = \int_{\Theta} P(\psi, T) P^{m-1}(\theta, d\psi), \quad m \in \mathbb{N}.$$

# Markovské procesy s obecnou množinou stavů

## Chapmanova-Kolmogorovova rovnost

### Věta 5.3 Chapmanova-Kolmogorovova rovnost.

V homogenním markovském řetězci s přechodovým jádrem  $P$  platí pro  $n, m \in \mathbb{N}_0$  a  $n \leq m$ ,  $\theta \in \Theta$ ,  $T \in \mathcal{T}$  vztah

$$P^m(\theta, T) = \int_{\Theta} P^{m-n}(\psi, T) P^n(\theta, d\psi).$$

*Důkaz.* V definici homogenního markovského řetězce stačí položit  $f_0 = \delta_{\theta}$ ,  $T_i = \Theta$  pro  $i = 0, \dots, m-1$ ,  $T_m = T$ . Definice  $P^n$  a  $P^{m-n}$  se použije pro prvních  $n$  a posledních  $m-n$  integrandů.



- Jak chápat Chapmanovu-Kolmogorovovu rovnost?

### Definice 5.6 Limitní rozdělení (*limitting distribution*).

Pravděpodobnostní rozdělení  $\pi$  na  $(\Theta, \mathcal{T})$  nazveme **limitní rozdělení** (*limitting distribution*) markovského řetězce  $\{\theta^{(m)} : m = 0, 1, \dots\}$  generovaného přechodovým jádrem  $P$ , jestliže

$$\lim_{m \rightarrow \infty} P^m(\theta, T) = \pi(T) \quad \text{pro } \pi\text{-s.v. } \theta \in \Theta \text{ a pro všechna } T \in \mathcal{T}.$$

- Poznámka. Je-li  $\pi$  limitním rozdělením, potom pro libovolné počáteční rozdělení  $f_0$  a pro každou  $T \in \mathcal{T}$  platí

$$P(\theta^{(m)} \in T) = \int_{\Theta} P^m(\theta, T) f_0(d\theta) \xrightarrow{m \rightarrow \infty} \int_{\Theta} \pi(T) f_0(d\theta) = \pi(T).$$

### Věta 5.4 Limitní a stacionární rozdělení.

*Je-li  $\pi$  limitním rozdělením homogenního markovského řetězce, potom je  $\pi$  též stacionárním rozdělením tohoto řetězce.*

*Důkaz.* Pro  $T \in \mathcal{T}$  a  $\pi$ -s.v.  $\theta \in \Theta$  je

$$\begin{aligned}\pi(T) &= \lim_{m \rightarrow \infty} P^m(\theta, T) = \lim_{m \rightarrow \infty} \int_{\Theta} P(\psi, T) P^{m-1}(\theta, d\psi) \\ &= \int_{\Theta} P(\psi, T) \pi(d\psi) = \pi P(T).\end{aligned}$$



## Oddíl **5.3**

# **Principy MCMC**

## Připomenutí, čím se zabýváme

- $f(d\theta)$  je nějaké pravděpodobnostní rozdělení.
- Pro měřitelné funkce  $t(\theta)$  potřebujeme aproximovat integrály typu

$$\int_{\Theta} t(\theta) f(d\theta) = \mathbb{E}_{f(d\theta)} t(\theta).$$

- Je-li  $S_M = \{\theta^{(1)}, \dots, \theta^{(M)}\}$  náhodný výběr z rozdělení  $f(d\theta)$ , potom (za jistých předpokladů)

$$\int_{\Theta} t(\theta) f(d\theta) \approx \frac{1}{M} \sum_{m=1}^M t(\theta^{(m)}) = \widehat{\mathbb{E}}_{f(d\theta)} t(\theta) := \widehat{t}_M.$$

▣ Monte Carlo integrace

- Necht'  $\{\theta^{(m)} : m = 0, \dots\}$  je homogenní markovský řetězec se **stacionárním** rozdělením  $f(d\theta)$ .
  - Víme: reversibilita vzhledem k  $f(d\theta)$  implikuje stacionaritu vzhledem k  $f(d\theta)$ .
- Stačí tedy zvolit přechodové jádro markovského řetězce tak, aby přechodová hustota  $p(\theta, \psi)$  splňovala **detailní podmínku rovnováhy** vzhledem k  $f(d\theta)$ .
- Stačí tedy volit přechodovou hustotu tak, aby splňovala

$$p(\theta, \psi) f(\theta) = p(\psi, \theta) f(\psi) \quad \forall \theta, \psi \in \Theta$$

---

a máme potřebný markovský řetězec.

▮▶ Toto není nikterak obtížné, jak bude záhy ukázáno.



- Předpokládejme, že se nám navíc podaří zajistit, že **existuje limitní** rozdělení uvažovaného markovského řetězce.
  - Víme: limitní rozdělení (existuje-li) = stacionární rozdělení  $f(d\theta)$ .
- Od jistého okamžiku (řekněme  $B + 1$ ) lze tedy
$$\mathcal{S}_M = \{\theta^{(B+1)}, \dots, \theta^{(B+M)}\}$$
považovat za náhodné veličiny s rozdělením  $f(d\theta)$ .
- Začátku řetězce  $\{\theta^{(0)}, \dots, \theta^{(B)}\}$  se říká **burn-in period**.
- Nejde o náhodný výběr, neboť  $\theta^{(B+1)}, \dots, \theta^{(B+M)}$  nejsou nezávislé!

- Nicméně, jestliže  $\int_{\Theta} |t(\theta)| f(d\theta) < \infty$  a jestliže dále platí **jisté předpoklady**, potom (**ergodická věta**):

$$\hat{t}_M = \frac{1}{M} \sum_{m=1}^M t(\theta^{(B+m)}) \xrightarrow{\text{a.s.}} \int_{\Theta} t(\theta) f(d\theta) \quad \text{pro } M \rightarrow \infty.$$

- $\hat{t}_M$  je tedy konzistentním odhadem pro  $\int_{\Theta} t(\theta) f(d\theta) = \mathbb{E}_{f(d\theta)} t(\theta)$ .
- Při splnění oněch **jistých předpokladů** lze též odhadnout

$$v_M = \text{var}_{f(d\theta)}(\hat{t}_M)$$

a odhadnout tak přesnost odhadu  $\mathbb{E}_{f(d\theta)} t(\theta)$

(přesnost aproximace integrálu  $\int_{\Theta} t(\theta) f(d\theta)$ ).

## Ergodická věta

- Předpoklady pro platnost ergodické věty pro markovské řetězce s obecnou množinou stavů jsou zobecněními předpokladů ergodické věty pro markovské řetězce s diskrétní množinou stavů.
- Potřeba zobecnit (a rozšířit) následující pojmy:
  - nerozložitelnost (*irreducibility*),
  - neperiodicita (*aperiodicity*),
  - trvalý (*recurrent*) a pozitivně trvalý (*positive recurrent*) markovský řetězec.
- ▣▶ **NMTP539: Metody Markov Chain Monte Carlo.**
- Zajistit splnění těchto předpokladů v praktických aplikacích též není těžké.
- Co je tedy obtížné?

### Největší obtíž při praktické aplikaci MCMC

- Zjistit, od kterého okamžiku již lze (s přiměřeně malou chybou) považovat rozdělení stavů vygenerovaného markovského řetězce za limitní = stacionární  $f(d\theta)$ .
  - Jak velké má být  $B$  (délka *burn-in period*)?
  - Jedná se o konvergenci pravděpodobnostních měr a nelze ji tedy posoudit jednoduchým číslem jako třeba v případě numerického optimalizačního algoritmu!

### Druhá největší obtíž při praktické aplikaci MCMC

- Připomeňme, že  $\theta^{(B+1)}, \dots, \theta^{(B+M)}$  nejsou obecně nezávislé a tudíž (i při předpokladu konvergence k limitnímu rozdělení) není nutně pravda

$$v_M = \text{var}_{f(d\theta)}(\hat{t}_M) = \frac{\text{var}_{f(d\theta)}(t(\theta))}{M}.$$

- Stavů markovského řetězce jsou typicky **kladně** (auto)korelovány a tudíž

$$v_M = \text{var}_{f(d\theta)}(\hat{t}_M) \geq \frac{\text{var}_{f(d\theta)}(t(\theta))}{M}.$$

- Je-li markovský řetězec zkonstruován tak, že vykazuje vysokou autokorelaci mezi jednotlivými stavy, může být  $v_M$  nepoužitelně vysoké i při poměrně vysoké hodnotě  $M$ .

## Největší obtíže

- Snaha konstruovat markovský řetězec tak, aby autokorelace byla co nejnižší.
  - Nulová autokorelace  $\equiv$  co do přesnosti se markovský řetězec chová stejně jako náhodný výběr, kde  $v_M = \frac{1}{M} \text{var}_{f(d\theta)}(t(\theta))$ .
- Konstrukce markovského řetězce s nízkou autokorelovaností snadná není a obtížnost této snahy závisí na konkrétní aplikaci.

## Oddíl **5.4**

# **Gibbsův algoritmus**

- Geman, S. and Geman, D. (1984). Stochastic relaxation, Gibbs distributions and the Bayes restoration of image. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **6**, 721–741.
  - Aplikace v oblasti restaurování digitálních obrázků.
- Gelfand, A. E. and Smith, A. F. M. (1990). Sampling-based approaches to calculating marginal densities. *Journal of the American Statistical Association*, **85**(410), 398–409.
  - Aplikace v bayesovské statistice.



## Předpoklady

### Předpoklady:

- $\Theta = \prod_{i=1}^k \Theta_i$ ,  $\theta = (\theta_1^\top, \dots, \theta_k^\top)^\top$
- Cílové (stacionární) rozdělení je  $f(d\theta)$  a má hustotu  $f(\theta)$  vzhledem k součinové míře  $\lambda_1 \otimes \dots \otimes \lambda_k$ , přičemž  $\lambda_i$  je  $\sigma$ -konečná míra s  $\lambda_i(\Theta_i) > 0$  ( $i = 1, \dots, k$ ).
- $\Theta = \{\theta : f(\theta) > 0\}$ .
- Jsme schopni (snadno) generovat z **plně podmíněných** (*full conditional*) rozdělení

$$f(d\theta_i | \theta_{-i}) = f(d\theta_i | \theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_k).$$

### Algoritmus:

1. Zvol počáteční stav  $\theta^{(0)} = (\theta_1^{(0)\top}, \dots, \theta_k^{(0)\top})^\top$ , polož  $m = 0$ .

2. (i) generuj  $\theta_1^{(m+1)}$  z podmíněného rozdělení  
 $f(d\theta_1 \mid \theta_2^{(m)}, \dots, \theta_k^{(m)})$ .

(ii) generuj  $\theta_2^{(m+1)}$  z podmíněného rozdělení  
 $f(d\theta_2 \mid \theta_1^{(m+1)}, \theta_3^{(m)}, \dots, \theta_k^{(m)})$ .

(iii) generuj  $\theta_3^{(m+1)}$  z podmíněného rozdělení  
 $f(d\theta_3 \mid \theta_1^{(m+1)}, \theta_2^{(m+1)}, \theta_4^{(m)}, \dots, \theta_k^{(m)})$ .

⋮

(k) generuj  $\theta_k^{(m+1)}$  z podmíněného rozdělení  
 $f(d\theta_k \mid \theta_1^{(m+1)}, \dots, \theta_{k-1}^{(m+1)})$ .

3. Zvětši  $m$  o jedničku a jdi na 2. krok algoritmu.

### Přechodová hustota

$$p(\theta, \psi) = \prod_{i=1}^k f(\psi_i | \psi_1, \dots, \psi_{i-1}, \theta_{i+1}, \dots, \theta_k).$$

- Odpovídá přechodovému jádru

$$P(\theta, T) = \int_T p(\theta, d\psi) = \int_T p(\theta, \psi) d\lambda(\psi).$$

### Věta 5.5 .

---

*Rozdělení  $f(d\theta)$  je stacionárním rozdělením markovského řetězce generovaného Gibbsovým algoritmem.*

---

*Důkaz.* Viz tabule.



- 
- Pokud bude existovat limitní rozdělení, musí se jednat o rozdělení stacionární a tedy cílové  $f(d\theta)$ .

## Existence limitního rozdělení, ergodicita

- **Ergodicitu** (existenci limitního rozdělení) lze dokázat například při splnění předpokladů, které byly uvedeny na začátku povídání o Gibbsově algoritmu, to jest
  - $\Theta = \prod_{i=1}^k \Theta_i$ ,  $\theta = (\theta_1^\top, \dots, \theta_k^\top)^\top$ .
  - Cílové (stacionární) rozdělení je  $f(d\theta)$  a má hustotu  $f(\theta)$  vzhledem k součinové míře  $\lambda_1 \otimes \dots \otimes \lambda_k$ , přičemž  $\lambda_j$  je  $\sigma$ -konečná míra s  $\lambda_j(\Theta_j) > 0$  ( $j = 1, \dots, k$ ).
  - $\Theta = \{\theta : f(\theta) > 0\}$ .
- Pro standardní statistické aplikace je toto obvykle splněno.
- Při rutinním použití Gibbsova algoritmu nicméně zůstává nemalým problémem zjistit, zda použitá **konečná** realizace markovského řetězce již dostatečně dobře odpovídá limitnímu = stacionárnímu = cílovému rozdělení.
- Při nevhodném použití Gibbsova algoritmu (viz dále) nemusí ani velmi dlouhá realizace markovského řetězce dostatečně dobře aproximovat limitní rozdělení!

## Reversibilita

- Lze ukázat, že markovský řetězec generovaný Gibbsovým algoritmem **ne-splňuje** detailní podmínku rovnováhy, tj. řetězec **není** reversibilní vzhledem k rozdělení  $f$ .
- Reversibility lze dosáhnout několika způsoby:
  - Generujeme střídavě **odpředu** a **odzadu**.
  - Pořadí vybíráme náhodně.
    - V každém podkroku Gibbsova algoritmu generujeme  $i$ -tý podvektor s pravděpodobnostmi  $p_i$  ( $0 < p_i < 1$ ,  $\sum_{i=1}^k p_i = 1$ ).
    - Častá volba je  $p_i = 1/k$  (rovnoměrné rozdělení).
    - **Random scan** Gibbsův algoritmus.

## Autokorelace

- V principu lze generovat ze všech **jednorozměrných** podmíněných rozdělení.
  - V případě, že složky  $\theta$  jsou v cílovém rozdělení  $f(d\theta)$  významně korelovány, vede generování z jednorozměrných podmíněných rozdělení k markovskému řetězci s velkou autokorelací.
  - Ideální situace je stav, kdy podvektory  $\theta_1, \dots, \theta_k$  jsou v cílovém rozdělení  $f(d\theta)$  co možná nejméně korelovány.

## Plně podmíněná rozdělení

- Při odvozování plně podmíněných rozdělení je vhodné si uvědomit a využívat základní fakt a to

$$f(\theta_i | \theta_{-i}) \propto f(\theta),$$

přičemž  $\propto$  nyní znamená, že vše, co neobsahuje  $\theta_i$  je konstantou.

- V případě **hierarchického** modelu, kde je  $f(\theta)$  zadáno jako součin postupně podmíněných rozdělení, pak  $f(\theta_i | \theta_{-i})$  závisí pouze na těch podmíněných rozděleních ze specifikace  $f(\theta)$ , kde jde v uvažované hierarchické struktuře:
  - o “potomky”  $\theta_i$ ,
  - o sourozence  $\theta_i$  (nejsou-li v  $f(d\theta)$  nezávislí),
  - o “rodiče”  $\theta_i$ .



## Lineární model

$$\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \mathbb{X} : \text{pevná matice } n \times k,$$
$$\boldsymbol{\varepsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n).$$

- Parametry:  $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \tau)^\top$ , kde  $\tau = \sigma^{-2} > 0$ .
- Věrohodnost:  $\mathbf{Y} \sim \mathcal{N}_n(\mathbb{X}\boldsymbol{\beta}, \tau^{-1} \mathbf{I}_n)$ .
- Neinformativní apriorní rozdělení:

$$p(\boldsymbol{\beta}) \propto 1, \quad \boldsymbol{\beta} \in \mathbb{R}^k,$$
$$p(\tau) \propto \frac{1}{\tau}, \quad \tau > 0.$$

## Příklad: Lineární model s neinformativním apriorním rozdělením

### Věrohodnost

- Označme:  $\mathbf{b} = (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbf{y}$ ,

$$SS_e = \|\mathbf{y} - \mathbb{X}\mathbf{b}\|^2.$$

- Věrohodnost:

$$L(\theta) = p(\mathbf{y} | \beta, \tau)$$

$$= (2\pi)^{-\frac{n}{2}} \tau^{\frac{n}{2}} \exp\left[-\frac{\tau}{2} \left\{ SS_e + (\beta - \mathbf{b})^\top \mathbb{X}^\top \mathbb{X} (\beta - \mathbf{b}) \right\}\right]$$

$$= (2\pi)^{-\frac{n}{2}} \tau^{\frac{n}{2}} \exp\left\{-\frac{\tau}{2} (\mathbf{y} - \mathbb{X}\beta)^\top (\mathbf{y} - \mathbb{X}\beta)\right\}.$$

# Příklad: Lineární model s neinformativním apriorním rozdělením

## Aposteriorní rozdělení

- Bylo odvozeno:

$$p(\boldsymbol{\beta}, \tau | \mathbf{y}) = p(\boldsymbol{\beta} | \tau, \mathbf{y}) \times p(\tau | \mathbf{y}),$$

kde  $p(\tau | \mathbf{y}) \sim \mathcal{G}\left(\frac{n-k}{2}, \frac{SS_e}{2}\right)$ ,

$$p(\boldsymbol{\beta} | \tau, \mathbf{y}) \sim \mathcal{N}_k\left(\mathbf{b}, \tau^{-1}(\mathbb{X}^T \mathbb{X})^{-1}\right).$$

- Dále bylo odvozeno:  $p(\boldsymbol{\beta} | \mathbf{y}) \sim \text{MVT}_{k, n-k}\left(\mathbf{b}, \frac{SS_e}{n-k}(\mathbb{X}^T \mathbb{X})^{-1}\right)$ .
- Pomocí Gibbsova algoritmu sestrojíme markovský řetězec, který bude mít rozdělení  $p(\boldsymbol{\beta}, \tau | \mathbf{y})$  jako stacionární i limitní.

## Příklad: Lineární model s neinformativním apriorním rozdělením

### Plně podmíněná rozdělení

- Označme  $\mathbb{W} = \mathbb{X}^\top \mathbb{X}$  s prvky  $w_{i,j}$  ( $i, j = 1, \dots, k$ ).

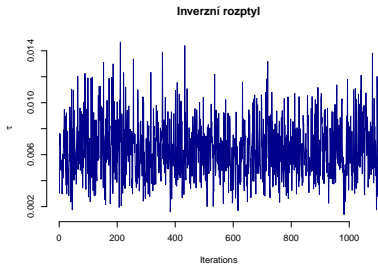
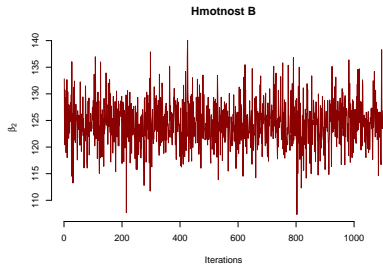
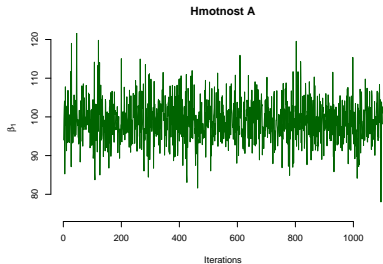
$$p(\boldsymbol{\beta} | \dots) = p(\boldsymbol{\beta} | \tau, \mathbf{y}) \sim \mathcal{N}_k(\mathbf{b}, \tau^{-1} (\mathbb{X}^\top \mathbb{X})^{-1}),$$

$$p(\beta_i | \dots) = p(\beta_i | \beta_{-i}, \tau, \mathbf{y}) \sim \mathcal{N}\left(b_i - \sum_{j \neq i} \frac{w_{i,j}}{w_{i,i}} (\beta_j - b_j), (\tau w_{i,i})^{-1}\right),$$

$$p(\tau | \dots) = p(\tau | \boldsymbol{\beta}, \mathbf{y}) \sim \mathcal{G}\left(\frac{n}{2}, \frac{(\mathbf{y} - \mathbb{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbb{X}\boldsymbol{\beta})}{2}\right).$$

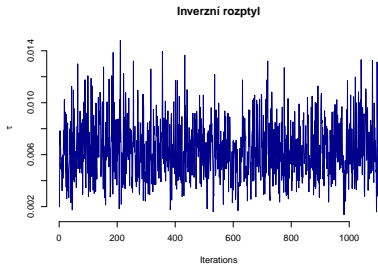
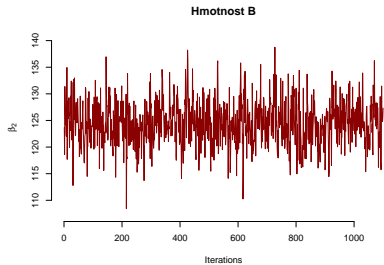
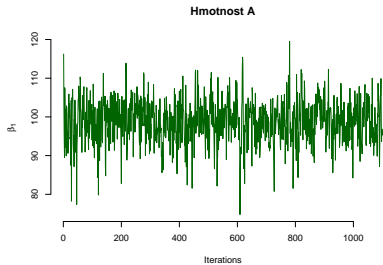
# Příklad: Vážení lehkých objektů

Blokový Gibbsův algoritmus: Generované hodnoty ( $B=100$ ,  $M=1\ 000$ )



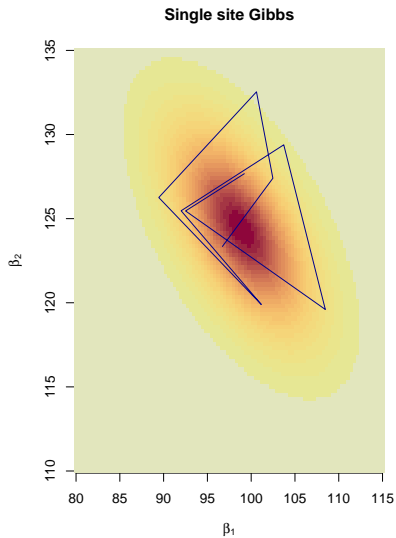
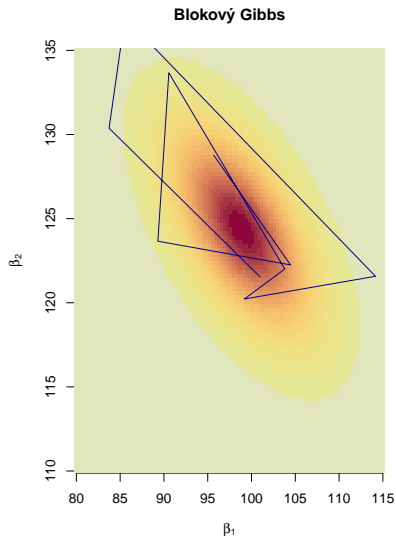
# Příklad: Vážení lehkých objektů

Single site Gibbsův algoritmus: Generované hodnoty ( $B=100$ ,  $M=1\ 000$ )



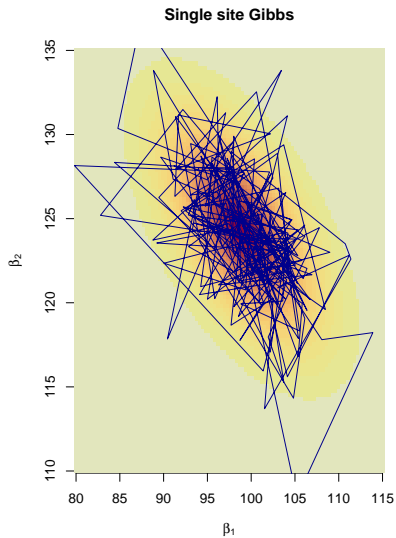
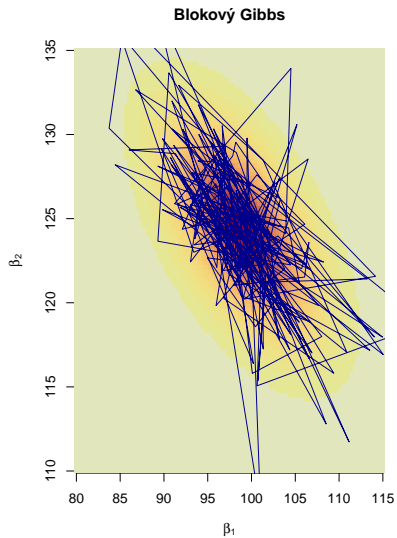
# Příklad: Vážení lehkých objektů

Gibbsův algoritmus: Generované hodnoty  $\beta$  (iterace 101 – 110)



# Příklad: Vážení lehkých objektů

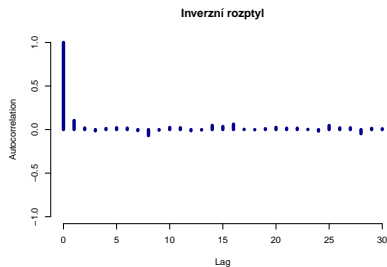
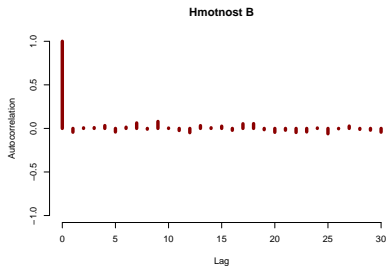
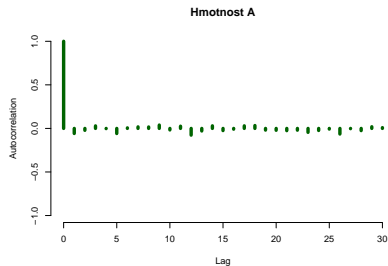
Gibbsův algoritmus: Generované hodnoty  $\beta$  (iterace 101 – 300)





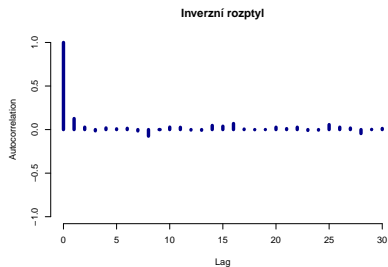
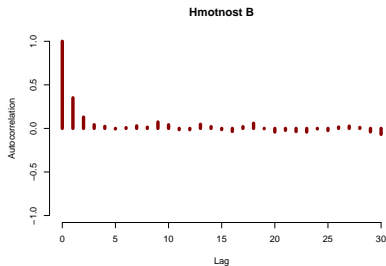
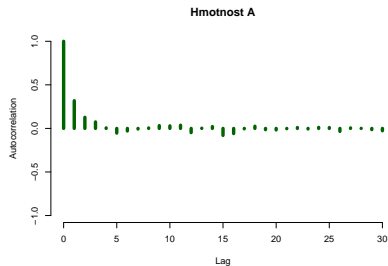
# Příklad: Vážení lehkých objektů

Blokový Gibbsův algoritmus: Odhady autokorelačních funkcí ( $B=100$ ,  $M=1\ 000$ )



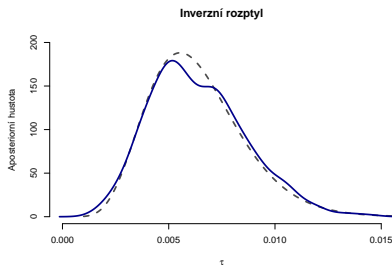
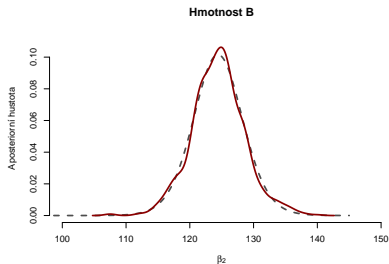
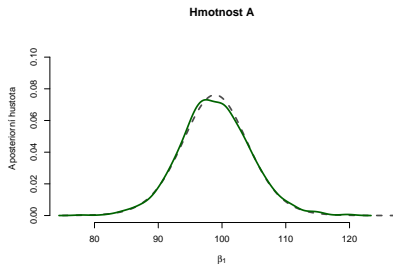
# Příklad: Vážení lehkých objektů

Single site Gibbsův algoritmus: Odhady autokorelačních funkcí ( $B=100$ ,  $M=1\ 000$ )



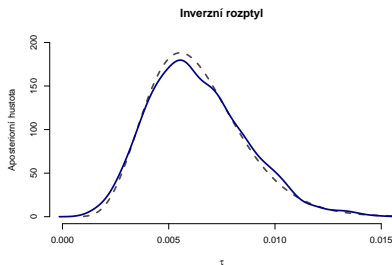
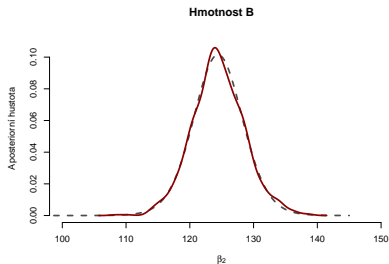
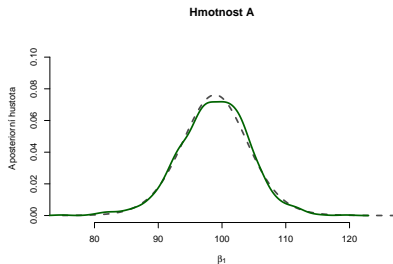
# Příklad: Vážení lehkých objektů

Blokový Gibbsův algoritmus: Odhady aposteriorních hustot ( $B=100$ ,  $M=1\ 000$ )



# Příklad: Vážení lehkých objektů

Single site Gibbsův algoritmus: Odhady aposteriorních hustot ( $B=100$ ,  $M=1\ 000$ )



# Příklad: Vážení lehkých objektů

Aposteriorní inference pro  $\beta$  (B=100, M=1 000)

## Blokový Gibbsův algoritmus

	$\beta_1$	$\beta_2$
<b>Aposter. střední hodnota</b>	98,8947	124,4211
<b>MCMC odhad</b>	98,9084	124,4561
<b>MC chyba (naivní)</b>	0,1744	0,1323
<b>MC chyba</b>	0,1662	0,1398
<b>Aposter. medián</b>	98,8947	124,4211
<b>MCMC odhad</b>	98,7209	124,4226
<b>95% ET věr. interval</b>	(87,9641; 109,8253)	(116,2231; 132,6190)
<b>MCMC odhad</b>	(86,9793; 110,2375)	(116,2702; 133,5293)
<b>95% HPD věr. interval</b>	(87,9641; 109,8253)	(116,2231; 132,6190)
<b>MCMC odhad</b>	(88,8421; 110,7950)	(114,8493; 132,0199)

# Příklad: Vážení lehkých objektů

Aposteriorní inference pro  $\beta$  (B=100, M=1 000)

## Single site Gibbsův algoritmus

	$\beta_1$	$\beta_2$
<b>Aposter. střední hodnota</b>	98,8947	124,4211
<b>MCMC odhad</b>	98,8051	124,4914
<b>MC chyba (naivní)</b>	0,1729	0,1302
<b>MC chyba</b>	0,2535	0,2015
<b>Aposter. medián</b>	98,8947	124,4211
<b>MCMC odhad</b>	98,9217	124,3193
<b>95% ET věr. interval</b>	(87,9641; 109,8253)	(116,2231; 132,6190)
<b>MCMC odhad</b>	(87,4650; 109,2495)	(116,3649; 133,4410)
<b>95% HPD věr. interval</b>	(87,9641; 109,8253)	(116,2231; 132,6190)
<b>MCMC odhad</b>	(87,8757; 109,4290)	(117,1579; 133,9186)

## Oddíl 5.5

# Metropolisův-Hastingsův algoritmus

- Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., and Teller, A. H. (1953). Equations of state calculations by fast computing machines. *Journal of Chemical Physics*, **21**, 1087–1091.
  - Aplikace ve statistické fyzice.
- Hastings, W. K. (1970). Monte Carlo sampling methods using Markov chains and their applications. *Biometrika*, **57**(1), 97–109.
  - Zobecnění algoritmu.
  - Uvážení též čistě statistických problémů.



## Předpoklady

### Předpoklady:

- Parametrický prostor  $\Theta$
- Cílové (stacionární) rozdělení je  $f(d\theta)$  a má hustotu  $f(\theta)$  vzhledem k  $\sigma$ -konečné míře  $\lambda$  s  $\lambda(\Theta) > 0$ .
- $\Theta = \{\theta : f(\theta) > 0\}$

# Metropolisův-Hastingsův algoritmus

## Algoritmus

### Algoritmus:

1. Zvol počáteční stav  $\theta^{(0)}$ , polož  $m = 0$ .
2. Generuj návrh  $\psi$  z rozdělení  $q(\theta^{(m)}, d\psi)$  s hustotou  $q(\theta^{(m)}, \psi)$  (vzhledem k  $\sigma$ -konečné míře  $\lambda$ ).
3. Spočti pravděpodobnost přijetí návrhu (*proposal acceptance probability*)

$$\alpha(\theta^{(m)}, \psi) = \begin{cases} \min \left\{ \frac{f(\psi) q(\psi, \theta^{(m)})}{f(\theta^{(m)}) q(\theta^{(m)}, \psi)}, 1 \right\} & \text{pro } f(\theta^{(m)}) q(\theta^{(m)}, \psi) > 0, \\ 1 & \text{jinak.} \end{cases}$$

4. Generuj  $U \sim \mathcal{U}(0, 1)$

$$\theta^{(m+1)} = \begin{cases} \psi, & \text{jestliže } U < \alpha(\theta^{(m)}, \psi), \\ \theta^{(m)}, & \text{jestliže } U \geq \alpha(\theta^{(m)}, \psi). \end{cases}$$

5. Zvětši  $m$  o jedničku a jdi na 2. krok algoritmu.

## Poznámky

### Poznámky:

- Pro aplikaci MH algoritmu není potřeba znát normující konstantu cílové hustoty  $f(\theta)$ .
  - Ideální pro použití v bayesovské statistice.
- Návrhová hustota  $q(\theta, \psi)$  může být **libovolná**.
  - Nevhodná volba  $q(\theta, \psi)$  však vede k vysoké autokorelaci a s tím spojené neefektivitě.
  - Příliš “ambiciózní”  $q(\theta, \psi)$  vede k malým pravděpodobnostem přijetí návrhu a řetězec pak dlouho setrvává v jednom stavu
    - ▣ vysoká autokorelace.
  - Příliš “opatrné”  $q(\theta, \psi)$  vede sice k vysokým pravděpodobnostem přijetí návrhu, ale řetězec se přesouvá jenom velice pomalu
    - ▣ vysoká autokorelace.
- Optimální proporce přijatých návrhů (*acceptance rate*) závisí na konkrétní situaci.

# Metropolisův-Hastingsův algoritmus

## Poznámky

- Symetrická návrhová hustota, tj.  $q(\theta, \psi) = q(\psi, \theta) \quad \forall \theta, \psi \in \Theta$ 
  - ▣▶ **Metropolisův** algoritmus.
- Hlavní část pravděpodobnosti přijetí

$$\alpha^*(\theta^{(m)}, \psi) = \frac{f(\psi) q(\psi, \theta^{(m)})}{f(\theta^{(m)}) q(\theta^{(m)}, \psi)}$$

obvykle počítáme v logaritmickém měřítku, tj.

$$\begin{aligned} \log\{\alpha^*(\theta^{(m)}, \psi)\} &= \log\{f(\psi)\} + \log\{q(\psi, \theta^{(m)})\} \\ &\quad - \log\{f(\theta^{(m)})\} - \log\{q(\theta^{(m)}, \psi)\} \end{aligned}$$

▣▶ vyhneme se mnoha numerickým obtížím při počítání s čísly, jež mohou být blízká nule.

## Návrhové rozdělení

Možné volby **návrhových rozdělení** (*proposal distribution*)

- **Nezávislý výběr** (*independent sampler*)

$$q(\theta, \psi) = q_0(\psi) \quad \forall \theta \in \Theta.$$

- 
- $q_0$  : nějaká hustota vzhledem k  $\sigma$ -konečné míře  $\lambda$  s nosičem na  $\Theta$ .
  - Návrhová hustota nezávisí na současném stavu.
  - Za  $q_0$  je vhodné volit rozdělení s těžšími chvosty (vícerozměrné t-rozdělení, ...).
  - Ideální stav:  $q_0(\psi) = f(\psi)$ 
    - generujeme přímo náhodný výběr z cílového rozdělení  $f(d\theta)$ .

## Návrhové rozdělení

Možné volby **návrhových rozdělení** (*proposal distribution*)

- **Náhodná procházka** (*random walk*)

$$q(\theta, \psi) = q_0(\psi - \theta) \quad \forall \theta \in \Theta.$$

- $q_0$  : nějaká hustota vzhledem k  $\sigma$ -konečné míře  $\lambda$  s nosičem na  $\Theta$ .
- Návrh:  $\psi = \theta + \mathbf{Z}$ , kde  $\mathbf{Z}$  má hustotu  $q_0$ .
- Častá volba:  $q_0 \equiv$  (vícerozměrné) normální, respektive t-rozdělení s nulovou střední hodnotou a obvykle diagonální varianční/měřtkovou maticí.
- ▣ Potřeba vhodně zvolit rozptyly.
  - Je-li  $q_0$  symetrická (tj.  $q_0(\mathbf{z}) = q_0(-\mathbf{z})$ ), potom  $q(\theta, \psi) = q(\psi, \theta)$  a při počítání pravděpodobnosti přijetí nemusíme vůbec počítat hodnoty hustoty  $q_0$  (resp. návrhové hustoty  $q$ ).

### Věta 5.6 .

---

*Rozdělení  $f(d\theta)$  je stacionárním rozdělením markovského řetězce generovaného Metropolisovým-Hastingsovým algoritmem.*

---

*Důkaz.* Viz přednáška NMTP539 Metody Markov Chain Monte Carlo.



- 
- Pokud bude existovat limitní rozdělení, musí se jednat o rozdělení stacionární a tedy cílové  $f(d\theta)$ .

# Metropolisův-Hastingsův algoritmus

---

## Existence limitního rozdělení, ergodicita

- Pro důkaz **ergodicity** (existence limitního rozdělení) je potřeba učinit několik předpokladů o návrhové hustotě  $q$ .
- Ergodicita je např. zajištěna v případě, kdy

$$q(\theta, \psi) = q_0(\psi - \theta), \quad q_0(\mathbf{z}) = q_0(-\mathbf{z})$$

---

(symetrická náhodná procházka)

a  $q_0(\mathbf{z}) > 0$  pro všechna  $\mathbb{R}^d$  (např. vícerozměrné normální nebo  $t$ -rozdělení).

- Podrobnosti viz přednáška NMTP539.



## Oddíl **5.6**

# Hybridní algoritmy

V Bayesovských aplikacích jsme obvykle v následující situaci.

- $\Theta = \prod_{i=1}^k \Theta_i$ ,  $\theta = (\theta_1^\top, \dots, \theta_k^\top)^\top$ .
- Cílové (stacionární) rozdělení je  $f(d\theta)$  a má hustotu  $f(\theta)$  vzhledem k součinové míře  $\lambda_1 \otimes \dots \otimes \lambda_k$ , přičemž  $\lambda_i$  je  $\sigma$ -konečná míra s  $\lambda_i(\Theta_i) > 0$  ( $i = 1, \dots, k$ ).
- $\Theta = \{\theta : f(\theta) > 0\}$ .
- Jsme schopni snadno vyjádřit všechna plně podmíněná rozdělení  $f(d\theta_i | \theta_{-i})$  až na multiplikativní konstantu, tj. víme, že pro hustoty platí

$$f(\theta_i | \theta_{-i}) \propto f^*(\theta_i | \theta_{-i}),$$

přičemž funkci  $f^*(\cdot | \theta_{-i})$  jsme schopni vyjádřit analyticky.

- Markovský řetězec se stacionárním rozdělením  $f(d\theta)$  lze sestavit tak, že jako základ vezmeme Gibbsův algoritmus a v případě, že pro nějaké  $i$  nejsme schopni (snadno) generovat z plně podmíněného rozdělení  $f(d\theta_i | \theta_{-i})$ , generujeme  $\theta_i$  pomocí Metropolisova-Hastingsova algoritmu
  - ▮ *Metropolis within Gibbs algorithm.*
- Celková přechodová hustota je **součinem** přechodových hustot pro Gibbsův, resp. Metropolisův-Hastingsův algoritmus.
- K tomu, aby celková procedura vedla k ergodickému řetězci s požadovaným limitním rozdělením  $f(d\theta)$  stačí, aby přechodová hustota v každém kroku vedla k ergodickému řetězci s požadovaným limitním rozdělením.

Některé další metody pro generování z plně podmíněných rozdělání

## ● Adaptive rejection sampling (ARS)

- Gilks, W. R. and Wild, P. (1992). Adaptive rejection sampling for Gibbs sampling. *Applied Statistics*, **41**, 337–348.
- Potřeba, aby hustota, z které chceme generovat byla **log-konkávní**.
- ▣ Poměrně častý případ, viz rozdělání z exponenciální třídy rozdělání.

## ● Adaptive rejection Metropolis sampling (ARMS)

- Gilks, W. R., Best, N. G., and Tan, K. K. C. (1995). Adaptive rejection Metropolis sampling within Gibbs sampling. *Applied Statistics*, **44**, 455–472.
- Zobecnění ARS metody na situace, kdy hustota, z které generujeme není log-konkávní.

## ● Slice sampling

- Neal, R. M. (2003). Slice sampling (with Discussion). *The Annals of Statistics*, **31**, 705–767.
- Efektivní v případě, že hustota, z které generujeme je **unimodální**.

## Lineární model

$$\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \mathbb{X} : \text{pevná matice } n \times k,$$
$$\boldsymbol{\varepsilon} \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n).$$

- Parametry:  $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \tau)^\top$ , kde  $\tau = \sigma^{-2} > 0$ .
- Věrohodnost:  $\mathbf{Y} \sim \mathcal{N}_n(\mathbb{X}\boldsymbol{\beta}, \tau^{-1} \mathbf{I}_n)$ .
- Neinformativní apriorní rozdělení:

$$p(\boldsymbol{\beta}) \propto 1, \quad \boldsymbol{\beta} \in \mathbb{R}^k,$$
$$p(\tau) \propto \frac{1}{\tau}, \quad \tau > 0.$$

## Příklad: Lineární model s neinformativním apriorním rozdělením

### Věrohodnost

- Označme:  $\mathbf{b} = (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbf{y}$ ,

$$SS_e = \|\mathbf{y} - \mathbb{X}\mathbf{b}\|^2.$$

- Věrohodnost:

$$L(\theta) = p(\mathbf{y} | \beta, \tau)$$

$$= (2\pi)^{-\frac{n}{2}} \tau^{\frac{n}{2}} \exp\left[-\frac{\tau}{2} \left\{ SS_e + (\beta - \mathbf{b})^\top \mathbb{X}^\top \mathbb{X} (\beta - \mathbf{b}) \right\}\right]$$

$$= (2\pi)^{-\frac{n}{2}} \tau^{\frac{n}{2}} \exp\left\{-\frac{\tau}{2} \|\mathbf{y} - \mathbb{X}\beta\|^2\right\}.$$

## Příklad: Lineární model s neinformativním apriorním rozdělením

### Apsteriorní rozdělení

- Bylo odvozeno:

$$p(\beta, \tau | \mathbf{y}) = p(\beta | \tau, \mathbf{y}) \times p(\tau | \mathbf{y}),$$

kde  $p(\tau | \mathbf{y}) \sim \mathcal{G}\left(\frac{n-k}{2}, \frac{SS_e}{2}\right)$ ,

$$p(\beta | \tau, \mathbf{y}) \sim \mathcal{N}_k\left(\mathbf{b}, \tau^{-1}(\mathbb{X}^T \mathbb{X})^{-1}\right).$$

- Dále bylo odvozeno:  $p(\beta | \mathbf{y}) \sim \text{MVT}_{k, n-k}\left(\mathbf{b}, \frac{SS_e}{n-k}(\mathbb{X}^T \mathbb{X})^{-1}\right)$
- Pomocí algoritmu Metropolis within Gibbs algoritmu sestrojíme markovský řetězec, který bude mít rozdělení  $p(\beta, \tau | \mathbf{y})$  jako stacionární i limitní.

## Příklad: Lineární model s neinformativním apriorním rozdělením

Plně podmíněná rozdělení

$$p(\beta | \dots) = p(\beta | \tau, \mathbf{y}) \sim \mathcal{N}_k(\mathbf{b}, \tau^{-1}(\mathbb{X}^\top \mathbb{X})^{-1}),$$

$$p(\tau | \dots) = p(\tau | \beta, \mathbf{y}) \sim \mathcal{G}\left(\frac{n}{2}, \frac{\|\mathbf{y} - \mathbb{X}\beta\|^2}{2}\right).$$



# Příklad: Lineární model s neinformativním apriorním rozdělením

## Metropolis within Gibbs algoritmus

- Regresní parametry  $\beta$  budeme generovat pomocí symetrické náhodné procházky s návrhovou hustotou

$$q(\beta_1, \beta_2) = q_0(\beta_2 - \beta_1),$$

kde  $q_0 \sim \mathcal{N}_k(\mathbf{0}, \mathbb{D}_{prop})$ ,  $\mathbb{D}_{prop} = \text{diag}(d_{1,prop}^2, \dots, d_{p,prop}^2)$ .

- V kroku  $m + 1$  algoritmu navrhneme  $\beta_{prop}$  vygenerované z rozdělení  $\mathcal{N}_k(\beta^{(m)}, \mathbb{D}_{prop})$ .
- Hlavní část pravděpodobnosti přijetí návrhu je

$$\begin{aligned} \alpha^*(\beta^{(m)}, \beta_{prop}) &= \frac{p(\beta_{prop} | \dots) q(\beta_{prop}, \beta^{(m)})}{p(\beta^{(m)} | \dots) q(\beta^{(m)}, \beta_{prop})} = \frac{p(\beta_{prop} | \dots)}{p(\beta^{(m)} | \dots)} \\ &= \exp \left[ -\frac{\tau^{(m)}}{2} \left\{ (\beta_{prop} - \mathbf{b})^\top \mathbb{X}^\top \mathbb{X} (\beta_{prop} - \mathbf{b}) - (\beta^{(m)} - \mathbf{b})^\top \mathbb{X}^\top \mathbb{X} (\beta^{(m)} - \mathbf{b}) \right\} \right]. \end{aligned}$$

## Příklad: Vážení lehkých objektů

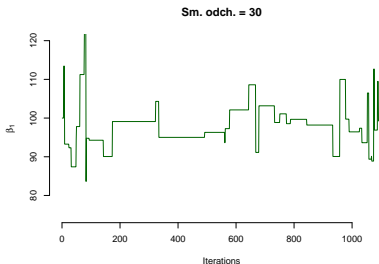
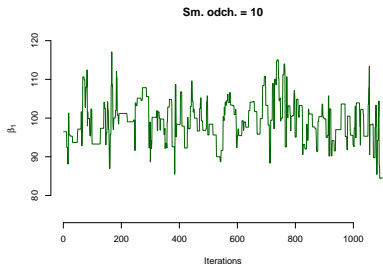
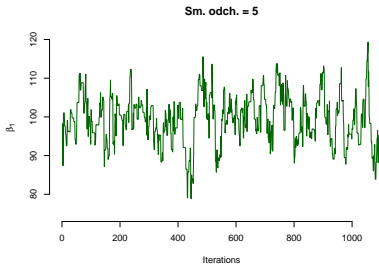
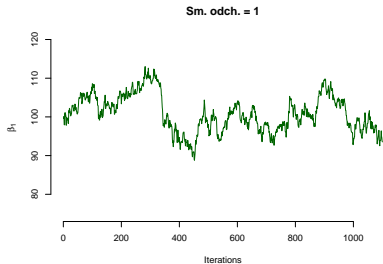
---

### Metropolis within Gibbs algoritmus

- Budou čtyři ukázky vygenerované při použití různých variančních matic v návrhové normální hustotě.
  1.  $\mathbb{D}_{prop} = \text{diag}(1, 1)$ 
    - ▣▶ proporce přijatých návrhů 0,87.
  2.  $\mathbb{D}_{prop} = \text{diag}(5^2, 5^2)$ 
    - ▣▶ proporce přijatých návrhů 0,47.
  3.  $\mathbb{D}_{prop} = \text{diag}(10^2, 10^2)$ 
    - ▣▶ proporce přijatých návrhů 0,22.
  4.  $\mathbb{D}_{prop} = \text{diag}(30^2, 30^2)$ 
    - ▣▶ proporce přijatých návrhů 0,04.

# Příklad: Vážení lehkých objektů

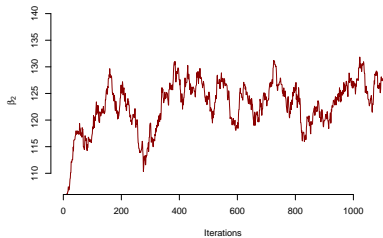
Metropolis within Gibbs algoritmus: Generované hodnoty  $\beta_1$  ( $B=100$ ,  $M=1\ 000$ )



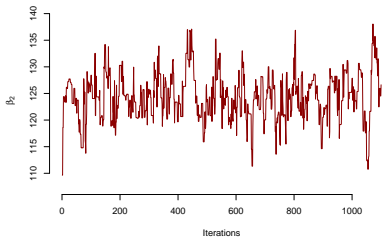
# Příklad: Vážení lehkých objektů

Metropolis within Gibbs algoritmus: Generované hodnoty  $\beta_2$  (B=100, M=1 000)

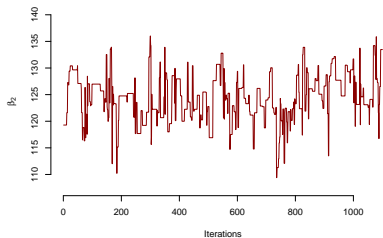
Sm. odch. = 1



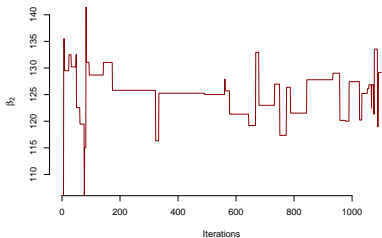
Sm. odch. = 5



Sm. odch. = 10

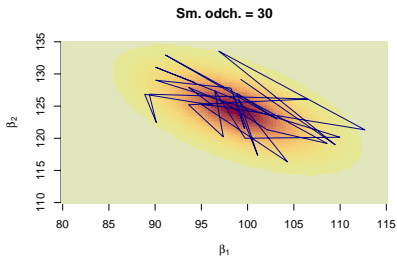
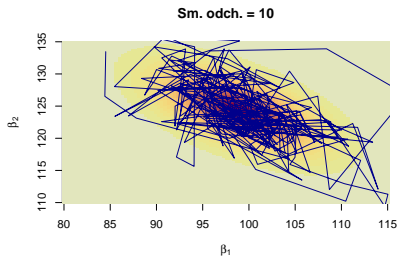
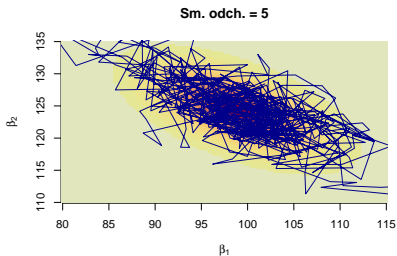
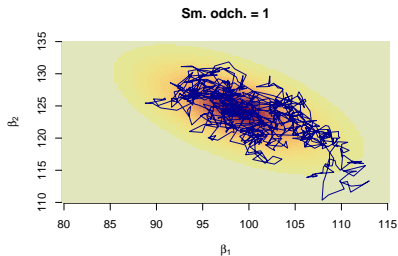


Sm. odch. = 30



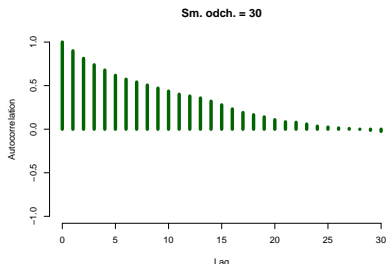
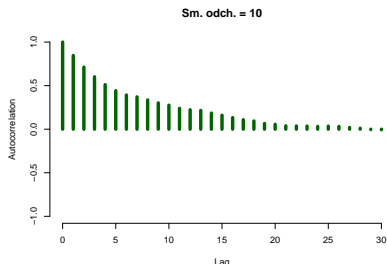
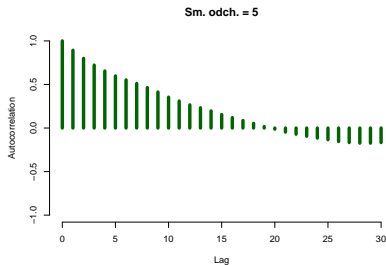
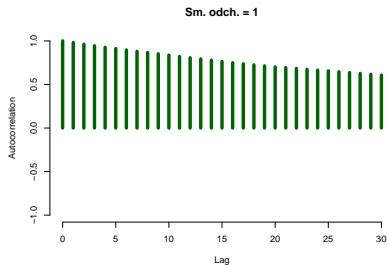
# Příklad: Vážení lehkých objektů

Metropolis within Gibbs algoritmus: Generované hodnoty  $\beta$  (iterace 101 – 1 100)



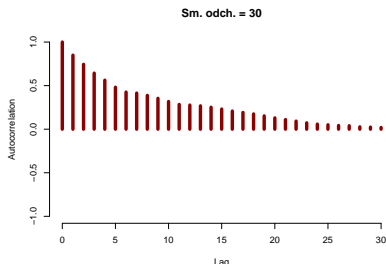
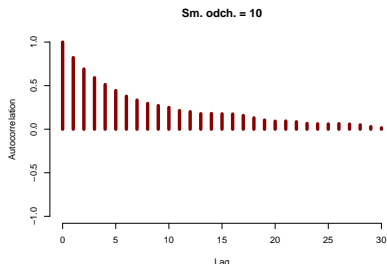
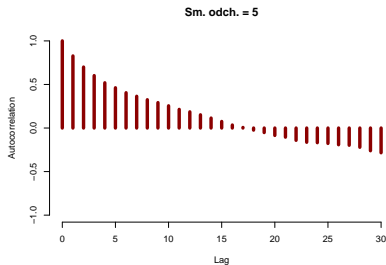
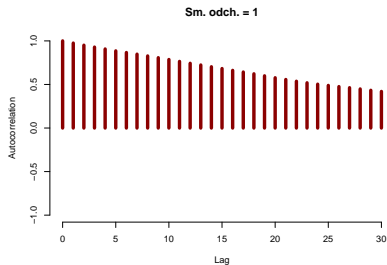
# Příklad: Vážení lehkých objektů

Metropolis within Gibbs algoritmus: Odhady autokorelačních funkcí pro  $\beta_1$  ( $B=100$ ,  $M=1\ 000$ )



# Příklad: Vážení lehkých objektů

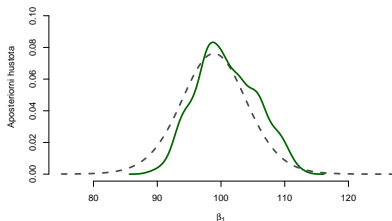
Metropolis within Gibbs algoritmus: Odhady autokorelačních funkcí pro  $\beta_2$  (B=100, M=1 000)



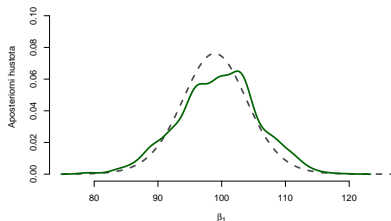
# Příklad: Vážení lehkých objektů

Metropolis within Gibbs algoritmus: Odhady aposteriorních hustot pro  $\beta_1$  (B=100, M=1 000)

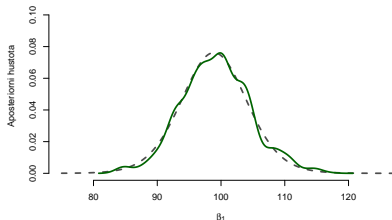
Smer. odch. = 1



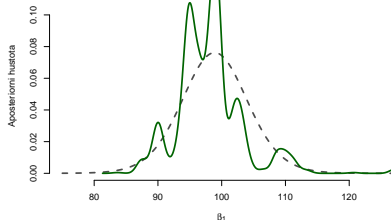
Smer. odch. = 5



Smer. odch. = 10



Smer. odch. = 30

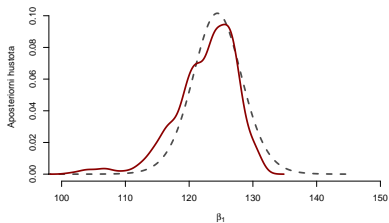




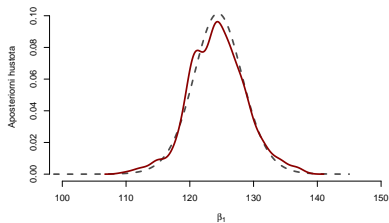
# Příklad: Vážení lehkých objektů

Metropolis within Gibbs algoritmus: Odhady aposteriorních hustot pro  $\beta_2$  ( $B=100$ ,  $M=1\ 000$ )

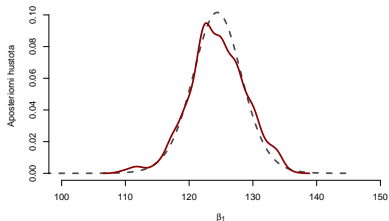
Smer. odch. = 1



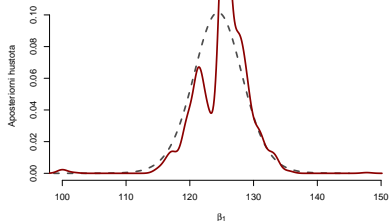
Smer. odch. = 5



Smer. odch. = 10



Smer. odch. = 30



# 6

## **Bayesovské rozšiřování dat**

# Oddíl **6.1**

## **Úvod**

## Rozšiřování dat (data augmentation)

- V bayesovské statistice potřebujeme typicky generovat z aposterioriního rozdělení s hustotou

$$p(\theta | \mathbf{y}) = \frac{L_{obs}(\theta) p(\theta)}{\int_{\Theta} L_{obs}(\theta) p(\theta) d\lambda(\theta)} \propto L_{obs}(\theta) p(\theta)$$

vzhledem k  $\sigma$ -konečné míře  $\lambda$  kde

- $L_{obs}(\theta) = p(\mathbf{y} | \theta)$  : věrohodnost (pozorovaných dat)
- $p(\theta)$  : apriorní rozdělení
- Pro použití MCMC metod obvykle:
  - není potřeba znát normující konstantu  $\int_{\Theta} L_{obs}(\theta) p(\theta) d\lambda(\theta)$ ;
  - je však vhodné a výhodné, aby výraz  $L_{obs}(\theta) p(\theta)$  byl snadno spočítatelný pro libovolné  $\theta \in \Theta$ .
- Často se však stává, že zejména  $L_{obs}(\theta)$  není jednoduchého tvaru.
  - nejčastější komplikace: při vyjadřování  $L_{obs}(\theta)$  je nutné integrovat.

# Příklad 1: Lineární smíšený model

**Model** (pro  $i = 1, \dots, N$ )

- $\mathbf{Y}_i | \mathbf{b}_i$  nezávislé s rozdělením  $\mathcal{N}_{n_i}(\mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \mathbf{b}_i, \sigma^2 \mathbf{I}_{n_i})$
- $\mathbf{b}_i$  i.i.d. s rozdělením  $\mathcal{N}_q(\boldsymbol{\mu}, \mathbb{D})$

**Parametry**

- $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \boldsymbol{\mu}^\top, \sigma^2, \text{vec}(\mathbb{D}))^\top$

**Věřohodnost** (pozorovaných dat)

$$\begin{aligned} L_{\text{obs}}(\boldsymbol{\theta}) &= p(\mathbf{y} | \boldsymbol{\theta}) = \prod_{i=1}^N p(\mathbf{y}_i | \boldsymbol{\theta}) = \prod_{i=1}^N \int_{\mathbb{R}^q} p(\mathbf{y}_i | \mathbf{b}_i, \boldsymbol{\theta}) p(\mathbf{b}_i | \boldsymbol{\theta}) d\mathbf{b}_i \\ &= \prod_{i=1}^N \int_{\mathbb{R}^q} \varphi(\mathbf{y}_i | \mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \mathbf{b}_i, \sigma^2 \mathbf{I}_{n_i}) \varphi(\mathbf{b}_i | \boldsymbol{\mu}, \mathbb{D}) d\mathbf{b}_i \end{aligned}$$

## Příklad 2: Logistická regrese s normálně rozdělenými náhodnými efekty

**Model** ( $j = 1, \dots, n_i$  pro každé  $i = 1, \dots, N$ )

- $Y_{i,j} | b_i$  nezávislé s rozdělením  $\mathcal{A}(\pi_i)$ , kde  $\pi_i = \frac{e^{b_i}}{1+e^{b_i}}$ .
- $b_i$  i.i.d. s rozdělením  $\mathcal{N}(\mu, d^2)$ .

**Parametry**

- $\theta = (\mu, d^2)^\top$ .

**Věrohodnost** (pozorovaných dat)

$$\begin{aligned} L_{obs}(\theta) &= p(\mathbf{y} | \theta) = \prod_{i=1}^N p(\mathbf{y}_i | \theta) = \prod_{i=1}^N \int_{\mathbb{R}} \prod_{j=1}^{n_i} p(y_{i,j} | b_i, \theta) p(b_i | \theta) db_i \\ &= \prod_{i=1}^N \int_{\mathbb{R}} \pi_i^{\sum_{j=1}^{n_i} y_{i,j}} (1 - \pi_i)^{n_i - \sum_{j=1}^{n_i} y_{i,j}} \varphi(b_i | \mu, d^2) db_i \\ &= \prod_{i=1}^N \int_{\mathbb{R}} \frac{e^{b_i \sum_{j=1}^{n_i} y_{i,j}}}{(1 + e^{b_i})^{n_i}} \varphi(b_i | \mu, d^2) db_i. \end{aligned}$$

## Příklad 3: Výsledky přijímacích zkoušek

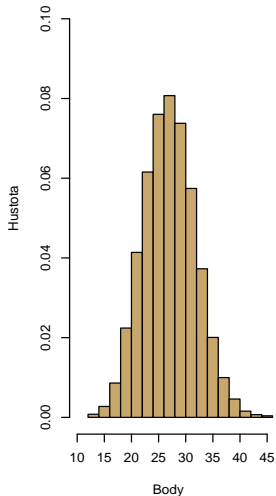
---

- V roce 1999 pořádala jistá (právnícká) fakulta jedné české VŠ (nebylo to v Plzni) 11 řádných a jeden náhradní termín (termín č. 12) přijímaček.
- Někdo si povšimnul, že u 12. náhradního termínu byla proporce přijatých studentů mnohem vyšší než u všech předchozích termínů.
  - Chytří studenti se hromadně omlouvali z řádného termínu přijímaček a přišli až na ten náhradní?
  - Výrazně více stimulující ovzduší v učebnách během 12. termínu?
  - Zázrak?
- Ukázalo, že zadání otázek pro 12. termín přijímaček záhadně uniklo (z uzamčeného trezoru) a bylo (v určitých kruzích) ke koupi již před konáním tohoto termínu.
- Celá událost byla dle tehdejšího rektora dané VŠ i děkana dotčené fakulty dílem *gangsterské mafie* stojící mimo fakultu.  
Viz <http://www.cibulka.com/nnoviny/nn2000/nn1900/obsah/05.htm>.

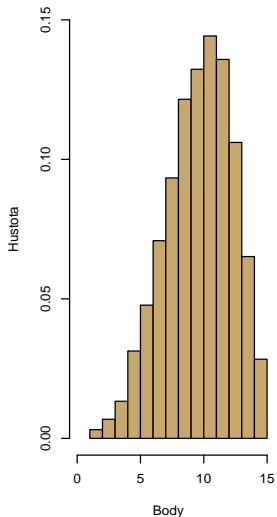
# Příklad 3: Výsledky přijímacích zkoušek

Termíny č. 1–11

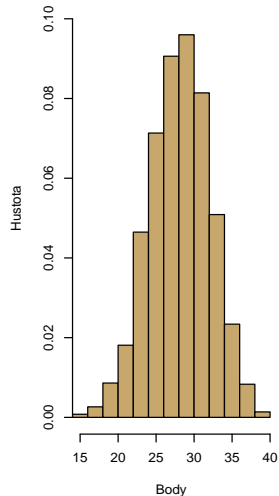
Historie



Jazyk



Logika

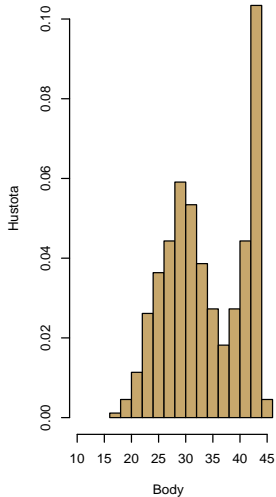




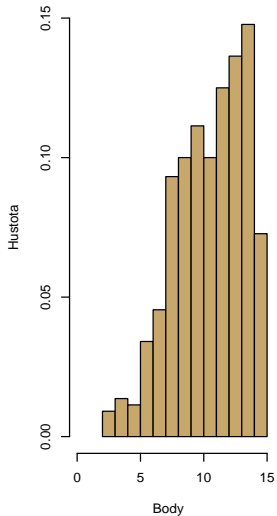
# Příklad 3: Výsledky přijímacích zkoušek

Termín č. 12

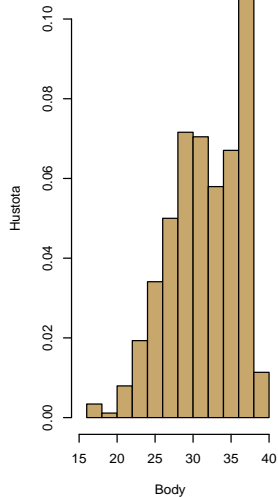
Historie



Jazyk



Logika



## Příklad 3: Výsledky přijímacích zkoušek

Možný model pro výsledky termínu č. 12

**Možný model** (pro  $i = 1, \dots, N$ )

- $\mathbf{Y}_i = (Y_{i,1}, Y_{i,2}, Y_{i,3})^\top$ : bodové zisky  $i$ tého studenta u jednotlivých částí zkoušky.
- Studenti pocházejí ze **dvou** populací:
  1. běžní studenti (proporce  $w_1$ );
  2. studenti napojení na gangsterskou mafii (proporce  $w_2$ ,  $w_1 + w_2 = 1$ ).
- Budeme **předpokládat**, že (sdružené) rozdělení bodových zisků je **v každé** populaci normální se střední hodnotou  $\mu_1$ , resp.  $\mu_2$  a varianční maticí  $\Sigma_1$ , resp.  $\Sigma_2$ .

⇒ Rozdělení  $\mathbf{Y}_i$ : **směs** normálních rozdělení s hustotou

$$p(\mathbf{y}_i | \theta) = w_1 \varphi(\mathbf{y}_i | \mu_1, \Sigma_1) + w_2 \varphi(\mathbf{y}_i | \mu_2, \Sigma_2)$$

## Příklad 3: Výsledky přijímacích zkoušek

Možný model pro výsledky termínu č. 12

- Potřeba odhadnout:
  - váhy (proporce)  $w_1, w_2$ ,
  - střední hodnoty  $\mu_1, \mu_2$ ,
  - varianční matice  $\Sigma_1, \Sigma_2$ .

$$\Rightarrow \theta = (w_1, w_2, \mu_1^\top, \mu_2^\top, \text{vec}(\Sigma_1), \text{vec}(\Sigma_2))^\top$$

**Věrohodnost** (pozorovaných dat)

$$L_{\text{obs}}(\theta) = p(\mathbf{y} | \theta) = \prod_{i=1}^N p(\mathbf{y}_i | \theta) = \prod_{i=1}^N \left\{ \sum_{k=1}^2 w_k \varphi(\mathbf{y}_i | \mu_k, \Sigma_k) \right\}.$$

## Oddíl **6.2**

# **Principy**

- Viděli jsme, že pozorovaná věrohodnost  $L_{obs}(\theta) = p(\mathbf{y} | \theta)$  (která tvoří základ aposteriorní hustoty  $p(\theta | \mathbf{y}) \propto L_{obs}(\theta) p(\theta)$ ) není vždy snadno vyjádřitelná jako **součin** “pěkných” funkcí.
- Někdy není  $L_{obs}(\theta)$  dokonce ani analyticky vyjádřitelná:
  - logistická regrese s náhodnými efekty,
  - zobecněné lineární smíšené modely (GLMM).
- Nicméně často se “věrohodnost” výrazně zjednoduší, jestliže budeme uvažovat více parametrů:
  - označme je  $\psi$ ;
  - “věrohodnost” je pak  $L_{augm}(\psi, \theta) = p(\mathbf{y} | \psi, \theta)$ ;
  - $L_{augm}(\psi, \theta)$  budeme nazývat **rozšířená** věrohodnost (*augmented likelihood*).

# Rozšiřování dat (data augmentation)

---

## Principy

- Parametry  $\psi$  mají často význam **nepozorovatelných** (resp. pouze **nepřímo pozorovatelných**) dat.
  - odsud termín **rozšiřování dat** (*data augmentation*);
  - termín pochází z článku Tanner, M. A. and Wong, W. H. (1987). The calculation of posterior distributions by data augmentation. *Journal of the American Statistical Association*, **82**(398), 528–550.
- $\mathbf{y}$  : pozorovaná/pozorovatelná data (*observed data*).
- $(\mathbf{y}, \psi)$  : úplná/rozšířená data (*complete data*).

# Rozšiřování dat (data augmentation)

## Principy

- Primárně nás zajímá  $p(\theta | \mathbf{y}) \propto \underbrace{p(\mathbf{y} | \theta)}_{L_{obs}(\theta)} p(\theta)$ ,

kde  $p(\mathbf{y} | \theta) = L_{obs}(\theta)$  plyne z předpokládaného (ne nutně hierarchického) modelu.

- Řekněme, že pro vhodné  $\psi$  se se sdruženou hustotou

$$p(\psi, \theta | \mathbf{y}) \propto p(\mathbf{y} | \psi, \theta) p(\psi, \theta) = \underbrace{p(\mathbf{y} | \psi, \theta)}_{L_{augm}(\psi, \theta)} p(\psi | \theta) p(\theta)$$

mnohem lépe pracuje.

- $L_{augm}(\psi, \theta)$ : model (věrohodnost) pro pozorovaná data, jestliže na doplněná data pohlížíme jako na další parametry modelu.
- $p(\psi | \theta)$ : model (věrohodnost) pro doplněná data.

## Principy

- Řekněme, že zajistíme, aby platilo

$$p(\theta | \mathbf{y}) = \int p(\psi, \theta | \mathbf{y}) d\lambda(\psi)$$

tj., aby  $p(\theta | \mathbf{y})$  bylo marginální hustotou rozdělení  $\theta | \mathbf{Y} = \mathbf{y}$  odpovídající sdruženému rozdělení  $(\psi, \theta) | \mathbf{Y} = \mathbf{y}$ .



## Principy

- Provádíme-li posteriorní inferenci na základě simulace, vygenerujeme náhodný výběr/markovský řetězec

$$\mathcal{S}_{(\psi, \theta), M} = \left\{ (\psi^{(1)}, \theta^{(1)}), \dots, (\psi^{(M)}, \theta^{(M)}) \right\}$$

s limitním rozdělením majícím hustotu  $p(\psi, \theta | \mathbf{y})$ .

- Jestliže  $p(\theta | \mathbf{y})$  je marginální hustotou odpovídající sdružené hustotě  $p(\psi, \theta | \mathbf{y})$ , potom

$$\mathcal{S}_{\theta, M} = \left\{ \theta^{(1)}, \dots, \theta^{(M)} \right\}$$

---

je náhodný výběr/markovský řetězec s limitním rozdělením majícím hustotu  $p(\theta | \mathbf{y})$ .

# Rozšiřování dat (data augmentation)

## Principy

- Máme:  $p(\theta | \mathbf{y}) \propto L_{obs}(\theta) p(\theta)$ ,  
 $p(\psi, \theta | \mathbf{y}) \propto L_{augm}(\psi, \theta) p(\psi | \theta) p(\theta)$ .
- $L_{obs}(\theta)$  plyne z předpokládaného modelu pro pozorovaná data.
- $L_{augm}(\psi, \theta) p(\psi | \theta)$  plyne z uvažovaného rozšiřování dat.
  - $L_{augm}(\psi, \theta)$ : model (věrohodnost) pro pozorovaná data, jestliže na doplněná data pohlížíme jako na další parametry modelu.
  - $p(\psi | \theta)$ : model (věrohodnost) pro doplněná data.
- Rozšiřování je potřeba udělat tak, aby  $p(\theta | \mathbf{y})$  bylo marginální hustotou odpovídající sdružené hustotě  $p(\psi, \theta | \mathbf{y})$ .
- Rozšiřování je tedy potřeba udělat tak, aby

$$L_{obs}(\theta) \propto \int L_{augm}(\psi, \theta) p(\psi | \theta) d\lambda(\psi).$$

# Rozšiřování dat (data augmentation)

## Principy

- K tomu, aby  $p(\boldsymbol{\theta} | \mathbf{y})$  bylo marginální hustotou odpovídající sdružené hustotě  $p(\boldsymbol{\psi}, \boldsymbol{\theta} | \mathbf{y})$  stačí, aby platilo

$$L_{obs}(\boldsymbol{\theta}) = p(\mathbf{y} | \boldsymbol{\theta}) \\ \propto \int p(\mathbf{y} | \boldsymbol{\psi}, \boldsymbol{\theta}) p(\boldsymbol{\psi} | \boldsymbol{\theta}) d\lambda(\boldsymbol{\psi}) = \int L_{augm}(\boldsymbol{\psi}, \boldsymbol{\theta}) p(\boldsymbol{\psi} | \boldsymbol{\theta}) d\lambda(\boldsymbol{\psi}).$$

- Výraz

$$p(\mathbf{y} | \boldsymbol{\psi}, \boldsymbol{\theta}) p(\boldsymbol{\psi} | \boldsymbol{\theta}) = L_{augm}(\boldsymbol{\psi}, \boldsymbol{\theta}) p(\boldsymbol{\psi} | \boldsymbol{\theta})$$

je roven  $p(\mathbf{y}, \boldsymbol{\psi} | \boldsymbol{\theta})$  a lze ho tedy interpretovat jako věrohodnost, jestliže bychom pozorovali **úplná** data.

- $L_{compl}(\boldsymbol{\theta}) := L_{augm}(\boldsymbol{\psi}, \boldsymbol{\theta}) p(\boldsymbol{\psi} | \boldsymbol{\theta})$   
= věrohodnost **úplných dat**.

# Rozšiřování dat (data augmentation)

## Principy

- K tomu, aby  $p(\theta | \mathbf{y})$  bylo marginální hustotou odpovídající sdružené hustotě  $p(\psi, \theta | \mathbf{y})$  tedy stačí specifikovat (rozšířený) model zahrnující nepozorovaná data  $\psi$  tak, aby si odpovídaly jednotlivé věrohodnosti.

► Přírozeně zajištěno v případě **hierarchických** modelů, kde rozšířená data  $\psi$  mají typicky přesně danou roli v popisu pravděpodobnostního mechanismu, o kterém předpokládáme, že generuje pozorovaná data  $\mathbf{y}$ .

Specifikace hierarchického modelu:

1.  $L_{augm}(\psi, \theta) = p(\mathbf{y} | \psi, \theta)$ : 1. hierarchická úroveň  
(model pro pozorovaná data za podmínky nepozorovaných dat).
2.  $p(\psi | \theta)$ : 2. hierarchická úroveň (model pro nepozorovaná data).
3.  $L_{obs}(\theta)$  (marginální model pro pozorovaná data)

“dopočítává” se jako  $L_{obs}(\theta) = \int L_{augm}(\psi, \theta) p(\psi | \theta) d\lambda(\psi)$ .

# Rozšiřování dat (data augmentation)

---

## Shrnutí terminologie

### Data, parametry

- $\mathbf{y}$  : pozorovaná data;
- $\theta$  : (frekventistické) parametry
  - inference o nich je naším primárním cílem;
- $\psi$  : nepozorovaná (pouze nepřímo pozorovaná) data, dodatečné parametry.

### Věrohodnosti

- $L_{obs}(\theta) = p(\mathbf{y} | \theta)$  : věrohodnost pozorovaných dat;
- $L_{augm}(\psi, \theta) = p(\mathbf{y} | \psi, \theta)$  : rozšířená věrohodnost pozorov. dat;
- $p(\psi | \theta)$  : věrohodnost doplněných dat;
- $L_{compl}(\theta) = p(\psi, \mathbf{y} | \theta) = L_{augm}(\psi, \theta) p(\psi | \theta)$  :  
věrohodnost úplných dat.

- Jednotlivé věrohodnosti potřeba specifikovat tak, aby

$$L_{obs}(\theta) = \int L_{augm}(\psi, \theta) p(\psi | \theta) d\lambda(\psi).$$

# Oddíl **6.3**

## **Příklady**

## Příklad 1: Lineární smíšený model

---

**Model** (pro  $i = 1, \dots, N$ )

- $Y_i \mid \mathbf{b}_i$  nezávislé s rozdělením  $\mathcal{N}_{n_i}(\mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \mathbf{b}_i, \sigma^2 \mathbf{I}_{n_i})$ ,
- $\mathbf{b}_i$  i.i.d. s rozdělením  $\mathcal{N}_q(\boldsymbol{\mu}, \mathbb{D})$ .

**Parametry**

- $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \boldsymbol{\mu}^\top, \sigma^2, \text{vec}(\mathbb{D}))^\top$ .

**Nepřímo pozorovatelná (doplněná) data**

- $\boldsymbol{\psi} = (\mathbf{b}_1^\top, \dots, \mathbf{b}_N^\top)^\top$ .

## Příklad 1: Lineární smíšený model

### Věrohodnost pozorovaných dat

$$\begin{aligned}L_{obs}(\boldsymbol{\theta}) &= \prod_{i=1}^N \int_{\mathbb{R}^q} \varphi(\mathbf{y}_i | \mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \mathbf{b}_i, \sigma^2 \mathbf{I}_{n_i}) \varphi(\mathbf{b}_i | \boldsymbol{\mu}, \mathbb{D}) d\mathbf{b}_i \\ &= \int_{\mathbb{R}^q} \cdots \int_{\mathbb{R}^q} \left\{ \prod_{i=1}^N \varphi(\mathbf{y}_i | \mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \mathbf{b}_i, \sigma^2 \mathbf{I}_{n_i}) \right\} \left\{ \prod_{i=1}^N \varphi(\mathbf{b}_i | \boldsymbol{\mu}, \mathbb{D}) \right\} d\mathbf{b}_1 \cdots d\mathbf{b}_N.\end{aligned}$$

### Rozšířená věrohodnost

$$L_{augm}(\boldsymbol{\psi}, \boldsymbol{\theta}) = \prod_{i=1}^N \varphi(\mathbf{y}_i | \mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \mathbf{b}_i, \sigma^2 \mathbf{I}_{n_i}).$$

### Věrohodnost doplněných dat

$$p(\boldsymbol{\psi} | \boldsymbol{\theta}) = \prod_{i=1}^N \varphi(\mathbf{b}_i | \boldsymbol{\mu}, \mathbb{D}).$$



## Příklad 2: Logistická regrese s normálně rozdělenými náhodnými efekty

---

**Model** ( $j = 1, \dots, n_i$  pro každé  $i = 1, \dots, N$ )

- $Y_{i,j} | b_i$  nezávislé s rozdělením  $\mathcal{A}(\pi_i)$ , kde  $\pi_i = \frac{e^{b_i}}{1+e^{b_i}}$ ,
- $b_i$  i.i.d. s rozdělením  $\mathcal{N}(\mu, d^2)$ .

**Parametry**

- $\theta = (\mu, d^2)^\top$ .

**Nepřímo pozorovatelná (doplněná) data**

- $\psi = (b_1, \dots, b_N)^\top$ .

## Příklad 2: Logistická regrese s normálně rozdělenými náhodnými efekty

### Věrohodnost pozorovaných dat

$$\begin{aligned}L_{obs}(\boldsymbol{\theta}) &= \prod_{i=1}^N \int_{\mathbb{R}} \frac{e^{b_i \sum_{j=1}^{n_i} y_{i,j}}}{(1 + e^{b_i})^{n_i}} \varphi(b_i | \mu, \sigma^2) db_i \\ &= \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} \left\{ \prod_{i=1}^N \frac{e^{b_i \sum_{j=1}^{n_i} y_{i,j}}}{(1 + e^{b_i})^{n_i}} \right\} \left\{ \prod_{i=1}^N \varphi(b_i | \mu, \sigma^2) \right\} db_1 \cdots db_N.\end{aligned}$$

### Rozšířená věrohodnost

$$L_{augm}(\boldsymbol{\psi}, \boldsymbol{\theta}) = \prod_{i=1}^N \frac{e^{b_i \sum_{j=1}^{n_i} y_{i,j}}}{(1 + e^{b_i})^{n_i}}.$$

### Věrohodnost doplněných dat

$$p(\boldsymbol{\psi} | \boldsymbol{\theta}) = \prod_{i=1}^N \varphi(b_i | \mu, \sigma^2).$$

## Příklad 3: Normální směšový model ( $K > 1$ skupin)

**Model** ( $i = 1, \dots, N$ )

- $Y_i$  i.i.d. se směšovým rozdělením s hustotou

$$p(\mathbf{y}_i | \theta) = \sum_{k=1}^K w_k \varphi(\mathbf{y}_i | \mu_k, \Sigma_k)$$

**Parametry**

- $\theta \equiv \mathbf{w}, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K$

$$\mathbf{w} = (w_1, \dots, w_K)^\top, 0 < w_k < 1, \sum_{k=1}^K w_k = 1$$

**Nepřímo pozorovatelná (doplněná) data**

- ???

### Věrohodnost pozorovaných dat

$$L_{\text{obs}}(\theta) = \prod_{i=1}^N \left\{ \sum_{k=1}^K \varphi(\mathbf{y}_i \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) w_k \right\}.$$

### Rozšířená věrohodnost

$$L_{\text{augm}}(\boldsymbol{\psi}, \theta) = \prod_{i=1}^N \varphi(\mathbf{y}_i \mid \boldsymbol{\mu}_{Z_i}, \boldsymbol{\Sigma}_{Z_i}).$$

### Věrohodnost doplněných dat

$$p(\boldsymbol{\psi} \mid \theta) = \prod_{i=1}^N w_{Z_i} = \prod_{i=1}^N \prod_{k=1}^K w_k^{\mathbb{I}(Z_i=k)}.$$

## Příklad 3: Normální směšový model

- Též nyní platí, že

$$L_{obs}(\theta) = \int L_{augm}(\psi, \theta) p(\psi | \theta) d\lambda(\psi).$$

- $\lambda$  je nyní součinnová číselná míra na  $\{1, \dots, K\}^N$  a tedy

$$\begin{aligned} & \int L_{augm}(\psi, \theta) p(\psi | \theta) d\lambda(\psi) \\ &= \sum_{z_1=1}^K \cdots \sum_{z_N=1}^K \left\{ \prod_{i=1}^N \varphi(\mathbf{y}_i | \boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i}) \right\} \left\{ \prod_{i=1}^N \underbrace{P(Z_i = z_i | \theta)}_{w_{z_i}} \right\} \\ &= \prod_{i=1}^N \left\{ \sum_{z_i=1}^K \varphi(\mathbf{y}_i | \boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i}) w_{z_i} \right\}. \end{aligned}$$

# Oddíl **6.4**

## **Poznámky**

# Rozšiřování dat (data augmentation)

## Poznámky

- Ani rozšiřování dat není úplně bez komplikací
- Zvětšujeme (často poměrně výrazně) **dimenzi** parametrického prostoru.
  - Generujeme-li z aposteriorního rozdělení pomocí MCMC, může být poměrně obtížné sestavit markovský řetězec s nízkou autokorelací a rychle konvergující k limitnímu rozdělení.
- Při použití rozšiřování dat pracujeme primárně s aposteriorním rozdělením

$$p(\psi, \theta | \mathbf{y}) \propto L_{augm}(\psi, \theta) p(\psi | \theta) p(\theta).$$

- Též zde lze apriorní rozdělení  $p(\theta)$  specifikovat hierarchicky za pomoci náhodných hyperparametrů  $\zeta$  s apriorní hustotou  $p(\zeta)$ .
- Fakticky potom pracujeme s aposteriorním rozdělením

$$p(\psi, \theta, \zeta | \mathbf{y}) \propto L_{augm}(\psi, \theta) p(\psi | \theta) p(\theta | \zeta) p(\zeta).$$

# Rozšiřování dat (data augmentation)

---

## Další oblasti využití

- Modely pro **cenzenovaná** data:
  - nejenom cenzenování zprava, ale též obecnější **intervalové** cenzenování,
  - nejenom neinformativní, ale též **informativní** cenzenování.
- A mnohé jiné. . .



# 7

## **Bayesian Model Selection**

## Oddíl 7.1

# Bayes factor

## Bayesian model

Data:  $\mathbf{Y}$ ,

Likelihood:  $p(\mathbf{y} | \boldsymbol{\theta}) = L(\boldsymbol{\theta}), \quad \boldsymbol{\theta} \in \Theta \subset \mathbb{R}^p,$

Prior distribution:  $p(\boldsymbol{\theta}).$

### Definice 7.1 Integrated (marginal) likelihood.

Marginal density of  $\mathbf{Y}$  following from the joint distribution of  $(\mathbf{Y}, \theta)$  is called the **integrated (marginal) likelihood**, i.e.,

$$p(\mathbf{y}) = \int_{\Theta} p(\mathbf{y}, \theta) d\theta = \int_{\Theta} \underbrace{p(\mathbf{y} | \theta)}_{L(\theta)} p(\theta) d\theta.$$

### Remarks

- Marginal likelihood is a likelihood of the model where the values of the unknown parameters are averaged over their prior distribution.
- It is also the denominator from the Bayes theorem.
- Also reported as **model evidence**.

# Model selection

---

- Interest in selecting a model from a set of candidate models  $M_1, \dots, M_r$ .

- Model  $M_k$ ,  $k = 1, \dots, r$ :

Likelihood:  $p_k(\mathbf{y} | \boldsymbol{\theta}_k) = L(\boldsymbol{\theta}_k), \quad \boldsymbol{\theta}_k \in \Theta_k \subset \mathbb{R}^{p_k},$

Prior distribution:  $p_k(\boldsymbol{\theta}_k),$

Integrated likelihood:  $p_k(\mathbf{y}) = \int_{\Theta} p_k(\mathbf{y} | \boldsymbol{\theta}_k) p_k(\boldsymbol{\theta}_k) d\boldsymbol{\theta}_k.$

- Integrated likelihood  $p_k(\mathbf{y})$  can also be interpreted as distribution of data under validity of model  $M_k$ :

$$p_k(\mathbf{y}) = p(\mathbf{y} | M_k), \quad k = 1, \dots, r.$$

## Model selection

- Let  $P(M_1), \dots, P(M_r)$  be the **prior** probabilities of models  $M_1, \dots, M_r$ :

$$0 < P(M_k) < 1, \quad k = 1, \dots, r \quad \sum_{k=1}^r P(M_k) = 1.$$

- For example (but not necessarily):  $P(M_k) = \frac{1}{r}, k = 1, \dots, r.$
- Model selection in Bayesian context can be based on **posterior** probabilities of models  $M_1, \dots, M_r$ :

$$P(M_k | \mathbf{y}) = \frac{p(\mathbf{y} | M_k) P(M_k)}{\sum_{l=1}^r p(\mathbf{y} | M_l) P(M_l)}, \quad k = 1, \dots, r.$$

- Choose model with the maximal posterior probability.
- “Small” complication: Integrated likelihood  $p_k(\mathbf{y}) = p(\mathbf{y} | M_k)$  must be calculated for each model which requires calculation of (usually complicated/intractable) integral.

## Definice 7.2 Bayes factor.

**Bayes factor** of the two models  $M_k$  and  $M_j$  is defined as the odds of the two integrated likelihoods, i.e.,

$$\text{BF}(M_k, M_j) = \frac{p(\mathbf{y} | M_k)}{p(\mathbf{y} | M_j)} = \frac{p_k(\mathbf{y})}{p_j(\mathbf{y})}.$$

### Remarks

- Bayes factor measures the evidence for model  $M_k$  versus model  $M_j$ .
- **Posterior odds** of model  $M_k$  versus model  $M_j$ :

$$= \frac{P(M_k | \mathbf{y})}{P(M_j | \mathbf{y})} = \frac{p(\mathbf{y} | M_k) P(M_k)}{p(\mathbf{y} | M_j) P(M_j)} = \text{BF}(M_k, M_j) \underbrace{\frac{P(M_k)}{P(M_j)}}_{\text{prior odds}(M_k, M_j)}.$$

- With the uniform prior distribution for the competing models:

$$\text{posterior odds}(M_k, M_j) = \text{BF}(M_k, M_j).$$

## Jeffreys' scale of evidence for Bayes factor

Bayes factor( $M_k, M_j$ )	Interpretation
$BF(M_k, M_j) < 1$	Negative support for $M_k$
$1 \leq BF(M_k, M_j) < 3$	Barely worth mentioning evidence for $M_k$
$3 \leq BF(M_k, M_j) < 10$	Substantial evidence for $M_k$
$10 \leq BF(M_k, M_j) < 30$	Strong evidence for $M_k$
$30 \leq BF(M_k, M_j) < 100$	Very strong evidence for $M_k$
$100 \leq BF(M_k, M_j)$	Decisive evidence for $M_k$



## Problems with Bayes factor

- The integrated likelihoods  $p_k(\mathbf{y})$  which enter the Bayes factor are, in fact, the **means (expected values)** of the likelihood (under model  $M_k$ ) with respect to the prior distribution (under model  $M_k$ ).
- $p_k(\mathbf{y})$  is not well defined when the prior distribution  $p_k(\theta_k)$  is **improper**.
- Bayes factor is numerically unstable when proper but **diffuse (weakly informative)** prior distributions used.
- There exist numerous approaches that were suggested in the literature to overcome above problems.

## Further reading

- Robert E. Kass, Adrian E. Raftery (1995). Bayes factors. *Journal of the American Statistical Association*. **90**(430), 773–795.
- Tomohiro Ando (2010). *Bayesian Model Selection and Statistical Modeling*. Boca Raton: Chapman & Hall/CRC. ISBN 978-1-4398-3614-9.

## Oddíl 7.2

# Posterior predictive distribution

## Bayesian model

Data:  $\mathbf{Y}$ ,

Likelihood:  $p(\mathbf{y} | \boldsymbol{\theta}) = L(\boldsymbol{\theta}), \quad \boldsymbol{\theta} \in \Theta \subset \mathbb{R}^p,$

Prior distribution:  $p(\boldsymbol{\theta}).$

## Posterior predictive distribution

- Let  $\mathbf{Y}_{new}$  be the random vector generated according to the same probabilistic mechanism as the data random vector  $\mathbf{Y}$ .
- In a Bayesian setting, it will always be assumed that  $\mathbf{Y}$  and  $\mathbf{Y}_{new}$  are (conditionally) independent given  $\theta$ .
- $\mathbf{Y}_{new} \equiv$  new (replicated) data.

### Definice 7.3 Posterior predictive distribution.

Posterior distribution of the random vector  $\mathbf{Y}_{new}$ , i.e.,  $p(\mathbf{y}_{new} | \mathbf{y})$ , is called the posterior predictive distribution.

We have

$$\begin{aligned} p(\mathbf{y}_{new} | \mathbf{y}) &= \int_{\Theta} p(\mathbf{y}_{new}, \theta | \mathbf{y}) d\theta = \int_{\Theta} p(\mathbf{y}_{new} | \theta, \mathbf{y}) p(\theta | \mathbf{y}) d\theta \\ &= \int_{\Theta} \underbrace{p(\mathbf{y}_{new} | \theta)}_{L_{new}(\theta)} p(\theta | \mathbf{y}) d\theta. \end{aligned}$$

## Integrated likelihood

$$p(\mathbf{y}) = \int_{\Theta} L(\theta) p(\theta) d\theta$$

- ≡ Distribution of data when the unknown parameters are averaged over their **prior** distribution.
- ▢ Evidence of the model **before** unknown parameters being estimated.

## Posterior predictive distribution

$$p(\mathbf{y}_{new} | \mathbf{y}) = \int_{\Theta} L_{new}(\theta) p(\theta | \mathbf{y}) d\theta$$

- ≡ Distribution of (new) data when the unknown parameters are averaged over their **posterior** distribution.
- ▢ Evidence of the model **after** using the data  $\mathbf{Y}$  for inference on unknown  $\theta$ .

## Oddíl 7.3

# Kullback-Leibler distance and deviance of the model

**Scetch of a theory will follow now which explains why the likelihood (or some of its derivatives) of the model can be considered as **evidence** of that model.**



### Definice 7.4 Kullback-Leibler distance.

Let  $Q_1$  and  $Q_2$  be two distributions with densities  $q_1$  and  $q_2$  (with respect to some  $\sigma$ -finite measure). The **Kullback-Leibler distance (divergence)** of  $Q_2$  from  $Q_1$  is defined as

$$\text{KL}(Q_2, Q_1) = \mathbb{E}_{Q_1} \log \left\{ \frac{q_1(\mathbf{Y})}{q_2(\mathbf{Y})} \right\} = \int q_1(\mathbf{y}) \log \left\{ \frac{q_1(\mathbf{y})}{q_2(\mathbf{y})} \right\} d\mathbf{y}.$$

- We have:  $\text{KL}(Q_2, Q_1) = \mathbb{E}_{Q_1} \log \{q_1(\mathbf{Y})\} - \mathbb{E}_{Q_1} \log \{q_2(\mathbf{Y})\}$ .
- Can also be shown:  $\text{KL}(Q_2, Q_1) \geq 0$ ,  
 $\text{KL}(Q_2, Q_2) = 0$ .

### In context of statistical modelling

- Let  $Q$  (with a density  $q$ ) be the **true** (unknown) distribution of data  $\mathbf{Y}$ .
- $L(\theta) = p(\cdot | \theta)$ : likelihood (**model**) for data (which possibly depends on a parameter vector  $\theta$ ).

Then

$$\text{KL}(L(\theta), Q) = \underbrace{\mathbb{E}_Q \log\{q(\mathbf{Y})\}}_{\text{const for all models}} - \mathbb{E}_Q \log\{p(\mathbf{Y} | \theta)\}.$$

- Up to an additive constant, the term  $-\mathbb{E}_Q \log\{p(\mathbf{Y} | \theta)\}$  is the Kullback-Leibler **distance** of the used model from the truth.

### Definice 7.5 Deviance of a model.

For given model with the likelihood  $L(\theta) = p(\mathbf{y} | \theta)$ , a quantity

$$D(\theta; \mathbf{y}) = -2 \log\{p(\mathbf{y} | \theta)\} = -2 \log\{L(\theta)\}$$

is called the **deviance** of the model.

### Remarks

- If  $Q$  is the true (unknown) distribution of data  $\mathbf{Y}$ , we have

$$2 \text{KL}(L(\theta), Q) = \mathbb{E}_Q\{D(\theta; \mathbf{Y})\} + \text{const.}$$

- Factor 2 in the definition of the deviance is used to get direct link to the statistic of the likelihood-ratio test.
- **Deviance**: suitable measure of the model quality (small deviance  $\equiv$  better model).

## Typically

Data:  $\mathbf{Y} = (\mathbf{Y}_1^\top, \dots, \mathbf{Y}_n^\top)^\top,$

Model (likelihood):  $L(\boldsymbol{\theta}) = p(\mathbf{y} | \boldsymbol{\theta}) = \prod_{i=1}^n p_i(\mathbf{y}_i | \boldsymbol{\theta}) = \prod_{i=1}^n L_i(\boldsymbol{\theta}),$

$\mathbf{Y}_1, \dots, \mathbf{Y}_n$  (conditionally) independent given  $\boldsymbol{\theta}$ .

## Deviance

$$\begin{aligned} D(\boldsymbol{\theta}; \mathbf{y}) &= -2 \log \left\{ \prod_{i=1}^n p_i(\mathbf{y}_i | \boldsymbol{\theta}) \right\} = -2 \sum_{i=1}^n \log \{ p_i(\mathbf{y}_i | \boldsymbol{\theta}) \} \\ &= \sum_{i=1}^n \underbrace{\left[ -2 \log \{ p_i(\mathbf{y}_i | \boldsymbol{\theta}) \} \right]}_{D_i(\boldsymbol{\theta}; \mathbf{y}_i)}. \end{aligned}$$

## Oddíl 7.4

# Measures of predictive ability of the model

**Our aim will now be to specify some criteria to evaluate/measure the ability of the model to make accurate **predictions** of new (replicated) data.**

**Those criteria can then be used for model selection.**

## Statistical model

(Observed) data:  $\mathbf{Y} = (\mathbf{Y}_1^\top, \dots, \mathbf{Y}_n^\top)^\top$ .

(New, not yet observed) data:  $\mathbf{Y}_{new} = (\mathbf{Y}_{new,1}^\top, \dots, \mathbf{Y}_{new,n}^\top)^\top$ .

$\mathbf{Y}$  and  $\mathbf{Y}_{new}$  generated by the same probabilistic mechanism.

---

Model (likelihood):  $L(\theta) = p(\cdot | \theta) = \prod_{i=1}^n p_i(\cdot | \theta)$ .

$\mathbf{Y}_1, \dots, \mathbf{Y}_n, \mathbf{Y}_{new,1}, \dots, \mathbf{Y}_{new,n}$  (conditionally) independent given  $\theta$ ,

$$p(\mathbf{y} | \theta) = \prod_{i=1}^n p_i(\mathbf{y}_i | \theta).$$

$$p(\mathbf{y}_{new} | \theta, \mathbf{y}) = p(\mathbf{y}_{new} | \theta) = \prod_{i=1}^n p_i(\mathbf{y}_{new,i} | \theta).$$

## Bayesian inference

Prior distribution:  $p(\theta)$ .

▣▶ Integrated likelihood:  $p(\mathbf{y}) = \int_{\Theta} p(\mathbf{y} | \theta) p(\theta) d\theta$

Evidence of the model before unknown parameters being estimated.

---

Inference on unknown  $\theta$  based on the posterior distribution:

$$p(\theta | \mathbf{y}) = \frac{p(\mathbf{y} | \theta) p(\theta)}{p(\mathbf{y})}.$$

▣▶ Posterior predictive distribution:  $p(\mathbf{y}_{new} | \mathbf{y}) = \int_{\Theta} p(\mathbf{y}_{new} | \theta) p(\theta | \mathbf{y}) d\theta$

Evidence of the model after the observed data  $\mathbf{Y} = \mathbf{y}$  used to infer on unknown parameters  $\theta$ .



## Posterior predictive deviance

Posterior predictive distribution:

$$p(\mathbf{y}_{new} | \mathbf{y}) = \int_{\Theta} p(\mathbf{y}_{new} | \theta) p(\theta | \mathbf{y}) d\theta = \int_{\Theta} \prod_{i=1}^n p_i(\mathbf{y}_{new,i} | \theta) p(\theta | \mathbf{y}) d\theta.$$

### Definition 7.6 Posterior predictive deviance.

Quantity

$$\bar{D}_{pred} = \mathbb{E}_{p(\theta | \mathbf{y})} D(\theta; \mathbf{y}_{new}) = \int_{\Theta} \underbrace{\left[ -2 \log \{ p(\mathbf{y}_{new} | \theta) \} \right]}_{D(\theta; \mathbf{y}_{new})} p(\theta | \mathbf{y}) d\theta$$

will be called the **posterior predictive deviance**.

- ▀ Suitable measure of prediction error (**loss of prediction**) when predicting  $\mathbf{Y}_{new} = \mathbf{y}_{new}$  using a (Bayesian) model estimated using data  $\mathbf{Y} = \mathbf{y}$ .

## Posterior predictive deviance

We have

$$\begin{aligned}\bar{D}_{pred} &= \mathbb{E}_{p(\theta | \mathbf{y})} D(\theta; \mathbf{y}_{new}) = \mathbb{E}_{p(\theta | \mathbf{y})} \left\{ \sum_{i=1}^n D_i(\theta; \mathbf{y}_{new,i}) \right\} \\ &= \sum_{i=1}^n \underbrace{\mathbb{E}_{p(\theta | \mathbf{y})} D_i(\theta; \mathbf{y}_{new,i})}_{\bar{D}_{pred,i}} = \sum_{i=1}^n \int_{\Theta} \underbrace{\left[ -2 \log \{ p_i(\mathbf{y}_{new,i} | \theta) \} \right]}_{D_i(\theta, \mathbf{y}_{new,i})} p(\theta | \mathbf{y}) d\theta.\end{aligned}$$

- To calculate  $\bar{D}_{pred}$  in practice (to be able to use it for model selection), we need the value of “new” data  $\mathbf{Y}_{new}$ .

### Posterior predictive deviance

$$\bar{D}_{pred} = \sum_{i=1}^n \bar{D}_{pred,i} = \sum_{i=1}^n \mathbb{E}_{p(\theta | \mathbf{y})} D_i(\theta, \mathbf{y}_{new,i}).$$

≡ Values of new data  $\mathbf{Y}_{new} = \mathbf{y}_{new}$  needed.

▣▶ Measure of the **loss of prediction**.

- Use **cross-validation** to estimate a value of each  $\bar{D}_{pred,i}$ ,  $i = 1, \dots, n$ :

$$\bar{D}_{pred,i} \approx \bar{D}_{pred,i}^{CV} = \mathbb{E}_{p(\theta | \mathbf{y}_{-i})} D_i(\theta, \mathbf{y}_i) = \int_{\Theta} D_i(\theta, \mathbf{y}_i) p(\theta | \mathbf{y}_{-i}) d\theta.$$

### Definice 7.7 Cross-validated posterior predictive deviance.

Quantity

$$\begin{aligned}\bar{D}_{pred}^{CV} &= \sum_{i=1}^n \underbrace{\mathbb{E}_{p(\theta | \mathbf{y}_{-i})} D_i(\theta; \mathbf{y}_i)}_{\bar{D}_{pred,i}^{CV}} \\ &= \sum_{i=1}^n \int_{\Theta} \underbrace{\left[ -2 \log \{ p_i(\mathbf{y}_i | \theta) \} \right]}_{D_i(\theta, \mathbf{y}_i)} p(\theta | \mathbf{y}_{-i}) d\theta.\end{aligned}$$

will be called the **cross-validated posterior predictive deviance**.

- With MCMC based Bayesian inference, (relatively) easily estimable if we have time to run the MCMC  $n$ -times (always with one observation left out).

### Definice 7.8 Posterior expected deviance.

Quantity

$$\bar{D} = \mathbb{E}_{p(\theta | \mathbf{y})} D(\theta; \mathbf{y}) = \int_{\Theta} \underbrace{\left[ -2 \log \{ p(\mathbf{y} | \theta) \} \right]}_{D(\theta; \mathbf{y})} p(\theta | \mathbf{y}) d\theta$$

will be called the **posterior expected deviance**.

We have

$$\begin{aligned} \bar{D} &= \mathbb{E}_{p(\theta | \mathbf{y})} D(\theta; \mathbf{y}) = \mathbb{E}_{p(\theta | \mathbf{y})} \sum_{i=1}^n D_i(\theta; \mathbf{y}_i) \\ &= \sum_{i=1}^n \underbrace{\mathbb{E}_{p(\theta | \mathbf{y})} D_i(\theta; \mathbf{y}_i)}_{\bar{D}_i} = \sum_{i=1}^n \int_{\Theta} \underbrace{\left[ -2 \log \{ p_i(\mathbf{y}_i | \theta) \} \right]}_{D_i(\theta, \mathbf{y}_i)} p(\theta | \mathbf{y}) d\theta. \end{aligned}$$

## Posterior expected deviance

$$\bar{D} = \sum_{i=1}^n \bar{D}_i = \sum_{i=1}^n \mathbb{E}_{p(\theta | \mathbf{y})} D_i(\theta, \mathbf{y}_i).$$

≡ Only the observed data  $\mathbf{Y} = \mathbf{y}$  needed.

▀ With MCMC based Bayesian inference, (relatively) easily estimable.

▀ Underestimates the **cross-validated posterior predictive deviance** which is

$$\bar{D}_{pred}^{CV} = \sum_{i=1}^n \bar{D}_{pred,i}^{CV} = \sum_{i=1}^n \mathbb{E}_{p(\theta | \mathbf{y}_{-i})} D_i(\theta, \mathbf{y}_i).$$

## Věta 7.1 .

For all  $i = 1, \dots, n$

$$\bar{D}_{pred,i}^{CV} - \bar{D}_i \geq 0.$$

## Reminder

$$\bar{D}_{pred,i}^{CV} = \mathbb{E}_{p(\theta | \mathbf{y}_{-i})} D_i(\theta, \mathbf{y}_i) = -2 \mathbb{E}_{p(\theta | \mathbf{y}_{-i})} \log p_i(\mathbf{y}_i | \theta),$$

$$\bar{D}_i = \mathbb{E}_{p(\theta | \mathbf{y})} D_i(\theta, \mathbf{y}_i) = -2 \mathbb{E}_{p(\theta | \mathbf{y})} \log p_i(\mathbf{y}_i | \theta).$$

$$\begin{aligned} \text{KL}_1 &:= \text{KL}\left(p(\boldsymbol{\theta} | \mathbf{y}), p(\boldsymbol{\theta} | \mathbf{y}_{-i})\right) \\ &= \int_{\Theta} p(\boldsymbol{\theta} | \mathbf{y}_{-i}) \log \left\{ \frac{p(\boldsymbol{\theta} | \mathbf{y}_{-i})}{p(\boldsymbol{\theta} | \mathbf{y})} \right\} d\boldsymbol{\theta} \\ &= \int_{\Theta} p(\boldsymbol{\theta} | \mathbf{y}_{-i}) \log \left\{ \frac{p(\mathbf{y}_{-i} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) p(\mathbf{y})}{p(\mathbf{y} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) p(\mathbf{y}_{-i})} \right\} d\boldsymbol{\theta} \\ &= \int_{\Theta} p(\boldsymbol{\theta} | \mathbf{y}_{-i}) \log \left\{ \frac{p(\mathbf{y})}{p_i(\mathbf{y}_i | \boldsymbol{\theta}) p(\mathbf{y}_{-i})} \right\} d\boldsymbol{\theta} \\ &= - \int_{\Theta} p(\boldsymbol{\theta} | \mathbf{y}_{-i}) \log \{ p_i(\mathbf{y}_i | \boldsymbol{\theta}) \} d\boldsymbol{\theta} + \log \left\{ \frac{p(\mathbf{y})}{p(\mathbf{y}_{-i})} \right\} \\ &= \frac{1}{2} \mathbb{E}_{p(\boldsymbol{\theta} | \mathbf{y}_{-i})} D_i(\boldsymbol{\theta}, \mathbf{y}_i) + \log \left\{ \frac{p(\mathbf{y})}{p(\mathbf{y}_{-i})} \right\} \\ &= \frac{1}{2} \overline{D}_{pred,i}^{CV} + \log \left\{ \frac{p(\mathbf{y})}{p(\mathbf{y}_{-i})} \right\}. \end{aligned}$$



$$\begin{aligned}\text{KL}_2 &:= \text{KL}\left(p(\boldsymbol{\theta} | \mathbf{y}_{-i}), p(\boldsymbol{\theta} | \mathbf{y})\right) \\ &= \int_{\Theta} p(\boldsymbol{\theta} | \mathbf{y}) \log \left\{ \frac{p(\boldsymbol{\theta} | \mathbf{y})}{p(\boldsymbol{\theta} | \mathbf{y}_{-i})} \right\} d\boldsymbol{\theta} \\ &= \int_{\Theta} p(\boldsymbol{\theta} | \mathbf{y}) \log \left\{ \frac{p(\mathbf{y} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) p(\mathbf{y}_{-i})}{p(\mathbf{y}_{-i} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) p(\mathbf{y})} \right\} d\boldsymbol{\theta} \\ &= \int_{\Theta} p(\boldsymbol{\theta} | \mathbf{y}) \log \left\{ \frac{p_i(\mathbf{y}_i | \boldsymbol{\theta}) p(\mathbf{y}_{-i})}{p(\mathbf{y})} \right\} d\boldsymbol{\theta} \\ &= \int_{\Theta} p(\boldsymbol{\theta} | \mathbf{y}) \log \{p_i(\mathbf{y}_i | \boldsymbol{\theta})\} d\boldsymbol{\theta} - \log \left\{ \frac{p(\mathbf{y})}{p(\mathbf{y}_{-i})} \right\} \\ &= -\frac{1}{2} \mathbb{E}_{p(\boldsymbol{\theta} | \mathbf{y})} D_i(\boldsymbol{\theta}, \mathbf{y}_i) - \log \left\{ \frac{p(\mathbf{y})}{p(\mathbf{y}_{-i})} \right\} \\ &= -\frac{1}{2} \bar{D}_i - \log \left\{ \frac{p(\mathbf{y})}{p(\mathbf{y}_{-i})} \right\}.\end{aligned}$$

That is,

$$\bar{D}_{pred,i}^{CV} - \bar{D}_i = 2 (\text{KL}_1 + \text{KL}_2) \geq 0.$$

### Definice 7.9 Expected optimism.

---

Quantity

$$\rho_{opt,i} = \mathbb{E}\left(\overline{D}_{pred,i}^{CV} - \overline{D}_i \mid \mathbf{Y}_{-i}\right), \quad i = 1, \dots, n$$

will be called the **expected optimism** when the loss of prediction of the  $i$ th observation is evaluated by  $\overline{D}_i$  ( $i$ th contribution to the posterior expected deviance) rather than by  $\overline{D}_{pred,i}^{CV}$  ( $i$ th cross-validated posterior predictive deviance).

### Definice 7.10 Penalized expected deviance (PED).

Quantity

$$\text{PED} = \underbrace{\sum_{i=1}^n \bar{D}_i}_{\bar{D}} + \underbrace{\sum_{i=1}^n \rho_{opt,i}}_{\rho_{opt}} = \sum_{i=1}^n \underbrace{(\bar{D}_i + \rho_{opt,i})}_{\text{PED}_i}$$

will be called the **penalized expected deviance (PED)**.

Quantity

$$\rho_{opt} = \sum_{i=1}^n \rho_{opt,i}$$

will be called the **overall optimism**, quantity

$$\text{PED}_i = \bar{D}_i + \rho_{opt,i}, \quad i = 1, \dots, n,$$

will be called contribution of the  $i$ th observation to the penalized expected deviance.

## Penalized expected deviance and cross-validated posterior predictive deviance

**Věta 7.2** Penalized expected deviance and cross-validated posterior predictive deviance.

The following holds for each  $i = 1, \dots, n$ :

$$\mathbb{E}(PED_i \mid \mathbf{Y}_{-i}) = \mathbb{E}(\overline{D}_{pred,i}^{CV} \mid \mathbf{Y}_{-i}).$$

With respect to cross-validation

$$PED = \sum_{i=1}^n PED_i$$

is equivalent to

$$\overline{D}_{pred}^{CV} = \sum_{i=1}^n \overline{D}_{pred,i}^{CV}.$$

## Penalized expected deviance and cross-validated posterior predictive deviance

---

*Důkaz.*

$$\begin{aligned}\mathbb{E}(\text{PED}_i \mid \mathbf{Y}_{-i}) &= \mathbb{E} \left\{ \bar{D}_i + \underbrace{\mathbb{E}(\bar{D}_{pred,i}^{CV} - \bar{D}_i \mid \mathbf{Y}_{-i})}_{\rho_{opt,i}} \mid \mathbf{Y}_{-i} \right\} \\ &= \mathbb{E}(\bar{D}_i \mid \mathbf{Y}_{-i}) + \mathbb{E}(\bar{D}_{pred,i}^{CV} \mid \mathbf{Y}_{-i}) - \mathbb{E}(\bar{D}_i \mid \mathbf{Y}_{-i}) \\ &= \mathbb{E}(\bar{D}_{pred,i}^{CV} \mid \mathbf{Y}_{-i}).\end{aligned}$$



- The last complication when using the PED for model comparison: calculation of the expected optimisms:

$$\rho_{opt,i} = \mathbb{E}\left(\bar{D}_{pred,i}^{CV} - \bar{D}_i \mid \mathbf{Y}_{-i}\right), \quad i = 1, \dots, n.$$

- 
- With MCMC based Bayesian inference, all expected optimisms  $\rho_{opt,i}$ ,  $i = 1, \dots, n$ , can be estimated using **two** parallel Markov chains (with  $p(\theta \mid \mathbf{y})$  as their limit distribution).

## Deviance information criterion

- For some classes of models, e.g., when  $p_i(\mathbf{y}_i | \theta)$ ,  $i = 1, \dots, n$ , belongs to **exponential** family, the overall optimism can be estimated as

$$p_{opt} = \sum_{i=1}^n p_{opt,i} \approx p_D = \bar{D} - D(\hat{\theta}(\mathbf{y}); \mathbf{y}),$$

where  $\hat{\theta}(\mathbf{y}) = \mathbb{E}_{p(\theta | \mathbf{y})} \theta$  (posterior mean of  $\theta$ ).

- Terminology:  $D(\hat{\theta}(\mathbf{y}); \mathbf{y})$ : **plug-in** deviance;  
 $p_D$ : effective number of parameters  
(measure of model complexity).
- “Small” inconvenience: the value of both the plug-in deviance and the effective number of parameters depends on the parameterization of the model.
- PED with  $p_D$  used in place of  $p_{opt}$   
▣ **Deviance information criterion (DIC).**



### Deviance information criterion (DIC)

$$\begin{aligned} \text{DIC} &= \bar{D} + p_D \\ &= \bar{D} + \left\{ \bar{D} - D(\hat{\theta}(\mathbf{y}); \mathbf{y}) \right\} \\ &= 2\bar{D} - D(\hat{\theta}(\mathbf{y}); \mathbf{y}). \end{aligned}$$

- DIC  $\equiv$  approximation to  $\bar{D}_{pred}^{CV}$  which was defined to evaluate the loss of prediction.
- Model with lower DIC is better.
- DIC is nowadays somehow overused/misused (even in situations when it is not justifiable)!

### Further reading

- David J. Spiegelhalter, Nicola G. Best, Bradley P. Carlin, Angelika Van Der Linde (2002). Bayesian measures of model complexity and fit (with Discussion). *Journal of the Royal Statistical Society, Series B*, **64**(4), 583–639.
- Martyn Plummer (2008). Penalized loss functions for Bayesian model comparison. *Biostatistics*, **9**(3), 523–539.
- David J. Spiegelhalter, Nicola G. Best, Bradley P. Carlin, Angelika Van Der Linde (2014). The deviance information criterion: 12 years on. *Journal of the Royal Statistical Society, Series B*, **76**(3), 485–493.