

# KOMPOZIČNÍ REGRESE S FUNKCIONÁLNÍ ZÁVISLE PROMĚNNOU

R. Talská, A. Menafoglio, K. Hron, J. Machalová, E. Fišerová



Katedra matematické analýzy a aplikací matematiky,  
Přírodovědecká fakulta, Univerzita Palackého v Olomouci

**11. - 16. 9. 2016**

- Motivace a cíl příspěvku
- Hustoty z pohledu Bayesových prostorů
- Funkcionální regresní model v  $L^2$  prostoru
- Funkcionální regresní model v Bayesově prostoru
- Shrnutí a výhled

- (Funkcionální) regresní analýza je široce užívána pro modelování **lineárního vztahu** mezi (funkcionální) **vysvětlovanou proměnnou a množinou** (skalárních) **proměnných**.
- **Otázka:** Je relevantní použít funkcionální modely navržené v  $L^2$  prostoru pro hustoty rozdělení pravděpodobností v roli vysvětlované proměnné?
- **Problém:**  $L^2$  prostor nezachycuje geometrické vlastnosti hustot.  
**Řešení:** Reprezentace hustot pomocí Bayesových prostorů.
- **Cíl:** Vytvořit funkcionální regresní model s využitím metodiky Bayesových prostorů.

- Hustoty jsou borelovsky měřitelné funkce (nezáporné, obvykle s jednotkovým integrálem), které obsahují informaci o relativních příspěvcích borelovských podmnožin nosiče hustot na celkovou míru tohoto nosiče.
- Bayesovy prostory  $\mathcal{B}^2$  obsahují ekvivalentní třídy hustot s nosičem na  $I = [a, b]$  (jejichž druhá mocnina logaritmu je integrovatelná).
- $\mathcal{B}^2$  má strukturu **separabilního Hilbertova prostoru**.
- **Centrovaná log-podílová transformace** (clr) = izometrický izomorfismus mezi  $\mathcal{B}^2$  a  $L^2$ . Pro  $f \in \mathcal{B}^2(I)$  máme:

$$f_c(t) = \text{clr}[f(t)] = \ln f(t) - \frac{1}{\eta} \int_I \ln f(t) dt, \quad t \in I.$$

# Funkcionální regresní model v $L^2$ prostoru

Předpokládejme vysvětlovanou proměnnou  $y(t)$  v  $L^2$  a  $p = r + 1$  nezávislých skalárních proměnných  $x_0, \dots, x_r$ .

$\mathbf{y}(t) = (y_1(t), \dots, y_N(t))$  ... vektor pozorovaných funkcí

$\mathbf{X} = [(x_{ij})]_{N \times p}$  ... matice plánu

Regrese s *absolutním členem*: první sloupec  $\mathbf{X}$  je tvořen 1.

**Funkcionální lineární model** je ve tvaru

$$y_i(t) = \beta_0(t) + \sum_{j=1}^r x_{ij} \beta_j(t) + \varepsilon_i(t), \quad i = 1, \dots, N, \quad (1)$$

kde  $\beta_j(t), j = 0, \dots, r$  jsou neznámé funkcionální regresní parametry v  $L^2$  a  $\varepsilon_i(t), i = 1, \dots, N$  jsou funkcionální náhodné chyby v  $L^2$ .

Model (1) maticově:  $\mathbf{y}(t) = \mathbf{X}\boldsymbol{\beta}(t) + \boldsymbol{\varepsilon}(t)$ .



# Funkcionální regresní model v $L^2$ prostoru

**Odhady**  $\hat{\beta}_j(t), j = 0, \dots, r$  získáme minimalizací následujících kriterií:

- 1 Reziduální součet čtverců (RSČ):

$$\text{RSČ}(\beta) = \int_I [\mathbf{y}(t) - \mathbf{X}\beta(t)]' [\mathbf{y}(t) - \mathbf{X}\beta(t)] dt, \quad (2)$$

- 2 Penalizované RSČ (PENRSČ):

$$\begin{aligned} \text{PENRSČ}(\beta) = & \int_I [\mathbf{y}(t) - \mathbf{X}\beta(t)]' [\mathbf{y}(t) - \mathbf{X}\beta(t)] dt + \\ & \lambda \int_I [L\beta(s)]' [L\beta(s)] ds. \end{aligned} \quad (3)$$

# Funkcionální regresní model v $\mathcal{B}^2$ prostoru

Předpokládejme vysvětlovanou proměnnou  $y(t)$  v  $\mathcal{B}^2$  a  $p = r + 1$  nezávislých skalárních proměnných  $x_0, \dots, x_r$ .

$\mathbf{y}(t) = (y_1(t), \dots, y_N(t))$  ... vektor pozorovaných hustot

$\mathbf{X} = [(x_{ij})]_{N \times p}$  ... matice plánu

Regrese s *absolutním členem*: první sloupec  $\mathbf{X}$  je tvořen 1.

**Funkcionální lineární model** je ve tvaru

$$y_i(t) = \beta_0(t) \oplus \bigoplus_{j=1}^r [x_{ij} \odot \beta_j(t)] \oplus \varepsilon_i(t), \quad i = 1, \dots, N, \quad (4)$$

kde  $\beta_j(t), j = 0, \dots, r$  jsou neznámé funkcionální regresní parametry v  $\mathcal{B}^2$  a  $\varepsilon_i(t), i = 1, \dots, N$  jsou funkcionální chyby nebo rezidua v  $\mathcal{B}^2$ .

# Funkcionální regresní model v $\mathcal{B}^2$ prostoru

**Odhady**  $\hat{\beta}_j(t), j = 0, \dots, r$  získáme použitím

$$\text{RSČ}(\beta) = \sum_{i=1}^N \|\varepsilon_i(t)\|_{\mathcal{B}}^2 = \sum_{i=1}^N \left\| \bigoplus_{j=0}^r [x_{ij} \odot \beta_j(t)] \ominus y_i(t) \right\|_{\mathcal{B}}^2 \quad (5)$$

Kritérium (5) v  $L^2$  prostoru:

$$\text{RSČ}(\beta) = \sum_{i=1}^N \|\text{clr}(\varepsilon_i(t))\|_2^2 = \sum_{i=1}^N \left\| \sum_{j=0}^r [x_{ij} \cdot \text{clr}(\beta_j(t))] - \text{clr}(y_i(t)) \right\|_2^2$$

s B-splajnovou reprezentací s omezující podmínkou na nulový integrál:

$$\text{clr}(y_i(t)) = \sum_k^K c_{ik} \varphi_k(t), \quad \int_1^K \sum_k^K c_{ik} \varphi_k(t) dt = 0, \quad (6)$$

$$\text{clr}(\beta_j(t)) = \sum_k^K b_{jk} \varphi_k(t) \quad \int_1^K \sum_k^K b_{jk} \varphi_k(t) dt = 0. \quad (7)$$



**Úloha:** nalézt **splajn s nulových integrálem** na  $I = [a, b]$ :

$$\int_a^b s_k(x) dx = 0.$$

→ podmínka je *přenesena* na podmínku pro B-splajnové koeficienty

**Věta (Nutná a postačující podmínka na  $\mathbf{b}$ )**

Pro splajn  $s_k(x) \in S_k^{\Delta\lambda}[a, b]$ ,  $s_k(x) = \sum_{i=-k}^g b_i B_i^{k+1}(x)$ ,

podmínka  $\int_a^b s_k(x) dx = 0$  je splněna tehdy a jen tehdy, když

$$\sum_{i=-k}^g b_i (\lambda_{i+k+1} - \lambda_i) = 0.$$

# Regrese s B-splajnovými koeficienty

**Vysvětlovaná**  $y_i(t)$ : splajny s koeficienty

$$\mathbf{Y}_{(i)} = (Y_{i,1}, \dots, Y_{i,g+k+1})' \text{ (řádky } \underline{\mathbf{Y}}_{N \times (g+k+1)})$$

**Mnohorozměrný lineární regresní model** je dán vztahem

$$\underline{\mathbf{Y}}_{(N \times (g+k+1))} = \mathbf{X}_{(N \times p)} \mathbf{B}_{(p \times (g+k+1))} + \underline{\boldsymbol{\varepsilon}}_{(N \times (g+k+1))} \quad (8)$$

Ekvivalentní zápis (8):

$$(\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_{g+k+1}) = \mathbf{X}(\beta_1, \beta_2, \dots, \beta_{g+k+1}) + (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_{g+k+1})$$

$\mathbf{Y}_{(i)} = (Y_{1,i}, \dots, Y_{g+k+1,i}), i = 1, \dots, N$  ... nezávislé  
mnohorozměrné proměnné se stejnou neznámou maticí  $\boldsymbol{\Sigma}$

$\mathbf{X}$  ... matice plánu s plnou sloupcovou hodnotí

$\underline{\boldsymbol{\varepsilon}}$  ... matice náhodných chyb

- Hustoty  $y_i(t)$ ,  $t \in I = [a, b]$  jsou reprezentovány pomocí kompozičních vektorů  $\mathbf{W}_i = (W_{i1}, \dots, W_{iD})'$ .
- Vyjádříme  $\mathbf{W}_i$  pomocí  $\mathbf{Z}_{(i)}$  v  $\mathbf{R}^D$  aplikováním clr transformace:

$$Z_{ij} = \ln \frac{W_{ij}}{\sqrt[D]{\prod_{j=1}^D W_{ij}}}, j = 1, \dots, D \rightarrow \underline{\mathbf{Z}}_{N \times D} = (\mathbf{Z}_{(1)}, \dots, \mathbf{Z}_{(N)})'$$

- Clr hustoty získáme vyhlazením clr dat  $\mathbf{Z}_{(i)}$ :

$$s_i(x) = \sum_{j=1}^{g+k+1} Y_{ij} B_j^{k+1}(x),$$

kde  $\mathbf{Y}_{(i)} = (Y_{i,1}, \dots, Y_{i,g+k+1})' = \mathbf{V}\mathbf{Z}_{(i)}$ ,  $i = 1, \dots, N$ ,  
nebo maticově:  $\underline{\mathbf{Y}}_{N \times (g+k+1)} = \underline{\mathbf{Z}}_{N \times D} \mathbf{V}'_{D \times (g+k+1)}$

- Predikované hodnoty:

$$\hat{\underline{Y}} = \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\underline{Y}, \quad \hat{\underline{Z}} = \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\underline{Z}$$

- Lze ukázat, že pokud predikujeme  $\underline{Y}$  (B-splajnových koeficientů), stejný výsledek získáme vyhlazením predikovaných clr dat  $\underline{Z}$ :

$$\begin{array}{ccc} \underline{Z} & \xrightarrow{\text{smoothing}} & \underline{Y} \\ \text{regression} \downarrow & & \downarrow \text{regression} \\ \hat{\underline{Z}} & \xrightarrow{\text{smoothing}} & \hat{\underline{Y}} \end{array}$$

→ kritérium PENRSČ je zahrnuto v RSČ kritériu díky použití vyhlazovacích splajnů.

- Použití přístupu Bayesových prostorů je nutnost v případě, že se hustota vyskytuje v roli vysvětlované proměnné ve funkcionálních lineárních regresních modelech.
- Odhad regresních funkcí (hustot) je možné převést na úlohu odhadu B-splajnových koeficientů  $c_l$  transformovaných hustot.  
(Funkcionální model se převede na mnohorozměrný model s B-splajnovými koeficienty v roli vysvětlované proměnné).
- Hladkost regresních odhadů (hustot)  $\hat{\beta}_j(t)$  je automaticky kontrolována vstupní B-splajnovou reprezentací.
- Výhled: aplikace teoretických poznatků na reálná data.

van den Boogaart, K.G., Egozcue, J.J., Pawlowsky-Glahn, V., *Bayes Hilbert spaces*. Australian & New Zeland Journal of Statistics 56(2), pp. 171-194, 2014.

Díaz-Barrero, J.L., Egozcue, J.J., Pawlowsky-Glahn, V., *Hilbert space of probability density functions based on Aitchison geometry*. Acta Mathematica Sinica, English Series 22, pp. 1175-1182, 2006.

Hron, K., Menafoglio, A., Templ, M., Hrušová, K., Filzmoser, P., *Simplicial principal component analysis for density functions in Bayes spaces*. Computational Statistics and Data Analysis 94, pp. 330-350, 2016.

Machalová, J., Hron, K., Monti, G.S., *Preprocessing of centred logratio transformed density functions using smoothing splines*. Journal of Applied Statistics, 2015, <http://dx.doi.org/10.1080/02664763.2015.11103706>.

Ramsay, J.O., Silverman, B.W., *Functional data analysis*. Springer, New York, 2005.