

# Kvantilová regrese

Oborový seminár - NMSA401

Lukáš Račko

4.1.2022



UNIVERZITA KARLOVA  
Matematicko-fyzikální  
fakulta

- 1 Lineárna regresia
  - Model
  - Problémy lineárneho modelu
- 2 Kvantilová regresia
  - Motivácia a model
  - Vlastnosti modelu
  - Výpočetné aspekty modelu
  - Inferencia o parametroch
- 3 Príklad
- 4 Zhrnutie

- $E(A|B)$  - Podmienená stredná hodnota  $A$  za podmienky  $B$ ;
- $\mathbf{Z}$  - Vektor vysvetľujúcich premenných;
- $Y$  - Závislá premenná;
- $\mathbb{X}$  - Regresná matica zodpovedajúca vhodnej transformácií vysvetľujúcich premenných;
- $\beta$  - Vektor parametrov modelu;
- $\hat{\theta}$  - Odhad parametru  $\theta$ ;
- $I_n$  - Matica identického zobrazenia na  $\mathbb{R}^n$ ;
- $\mathbf{1}_n$  - Vektor jednotiek dimenzie  $n$ ;
- $\|\cdot\|_2$  -  $\ell^2$ -norma vektoru.

Nech  $\begin{pmatrix} Y_1 \\ \mathbf{Z}_1 \end{pmatrix}, \dots, \begin{pmatrix} Y_n \\ \mathbf{Z}_n \end{pmatrix}$  je náhodný výber. Označíme  $\mathbf{Y} = (Y_1, \dots, Y_n)$ :

- Predpokladáme  $\mathbf{Y}|\mathbb{X} \sim (\mathbb{X}\boldsymbol{\beta}, \sigma^2 I_n)$ ;
- Modelujeme  $E(Y|\mathbf{Z})$ ;
- Riešime  $\arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \|\mathbf{Y} - \hat{\mathbf{Y}}(\boldsymbol{\beta})\|_2 = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \|\mathbf{Y} - \mathbb{X}\boldsymbol{\beta}\|_2$ .

Príliš silné predpoklady modelu.

- Normalita - väčšinou nie až tak dôležitá;
- Homoskedasticita.

Jednoduchý model vytvára problémy pri interpretácii zložitejších dát a môže byť absolútne nevhodný pre riešenie niektorých problémov.

- Riešenie heteroskedasticity - Box-Coxové transformácie;
- Zanedbáva "riziková štruktúra" problému.

Chceme vytvoriť model, ktorý bude modelovať distribúciu  $Y|Z$ . Dôvody:

- Vyriešime tak problém s interpretáciou heteroskedastických dát;
- Budeme "poznať" riziko v probléme. Dôležité pri modelovaní vo financiách.

- Pre  $\tau \in (0, 1)$  definujeme stratovú funkciu

$$\rho_\tau(x) = \tau x \mathbb{I}\{x > 0\} - (1 - \tau)x \mathbb{I}\{x \leq 0\}$$

- a jej "deriváciu" ako

$$\psi_\tau(x) = \rho'_\tau(x) = \tau \mathbb{I}\{x > 0\} - (1 - \tau) \mathbb{I}\{x < 0\}, x \neq 0$$

a  $\psi_\tau(0) := 0$ .

Nech náhodná veličina  $X : (\Omega, \mathcal{A}) \rightarrow (\mathbb{R}, \mathcal{B})$  má distribučnú funkciu  $F$ , potom

$$\forall \tau \in (0, 1) : F^{-1}(\tau) = \arg \min_{\theta \in \mathbb{R}} E [\rho_{\tau}(X - \theta) - \rho_{\tau}(X)].$$

Pre náhodný výber  $X_1, \dots, X_n \sim X$  potom zjavne platí

$$\forall \tau \in (0, 1) : \hat{F}_n^{-1}(\tau) = \arg \min_{\theta \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \rho_{\tau}(X_i - \theta).$$

Majme náhodný výber

$$\begin{pmatrix} \mathbf{Z}_1 \\ Y_1 \end{pmatrix}, \dots, \begin{pmatrix} \mathbf{Z}_n \\ Y_n \end{pmatrix}$$

z nejakého rozdelenia náhodného vektoru  $\begin{pmatrix} \mathbf{Z} \\ Y \end{pmatrix}$ . Pre  $\tau \in (0, 1)$  definujeme  $\tau$ -tý regresný kvantil ako

$$\beta_{\mathbf{Z}}(\tau) = \arg \min_{\mathbf{b} \in \mathbb{R}^p} E(\rho_{\tau}(Y - \mathbf{X}^T \mathbf{b}) - \rho_{\tau}(Y))$$

a jeho odhad ako

$$\hat{\beta}_n(\tau) = \arg \min_{\mathbf{b} \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \rho_{\tau}(Y_i - \mathbf{X}_i^T \mathbf{b})$$

Tak ako v lineárnom modele, opäť predpokladáme lineárnu závislosť modelovanej kvantily na regresoroch. T.j. predpokladáme

$$F_{Y|Z}^{-1}(\tau) = \mathbf{X}^T \beta_{\mathbf{X}}(\tau).$$

V lineárnom modele sme predpokladali

$$E(Y|Z) = \mathbf{X}^T \beta.$$

Za platnosti modelu  $F_{Y|Z}^{-1}(\tau) = \mathbf{X}^T \boldsymbol{\beta}$  máme "štandardnú" interpretáciu hodnôt  $\hat{\beta}_{nk}(\tau)$  ako odhad zmeny hodnoty  $\tau$ -kvantilu podmienenej distribúcie  $Y|Z$  pri zmene  $k$ -tého regresoru o 1. Špeciálne v prípade jednoduchého modelu s interceptom to znamená zmenu  $k$ -tej vysvetľujúcej premennej o 1.

Špeciálne sa pozrime na prípad  $\tau = 0.5$ . Vtedy dostávame  $\rho_\tau(\cdot) = 0.5|\cdot|$  a preto zjavne

$$\beta_{\mathbf{X}} = \arg \min_{\mathbf{b} \in \mathbb{R}^p} E (|Y - \mathbf{X}^T \mathbf{b}| - |Y|).$$

Takže  $\mathbb{L}_1$ -regresia je špeciálnym prípadom kvantilovej regresie a zodpovedá regresnému mediánu.

Uvažujeme homoskedastický model s interceptom

$$Y \sim \beta_0 + \mathbf{X}^T \boldsymbol{\beta} + \varepsilon,$$

kde  $\varepsilon \perp \mathbf{X}$ . Teda pre  $\tau \in (0, 1)$  pevné platí

$$F_{Y|Z}^{-1}(\tau) = \beta_0 + F_{\varepsilon}^{-1}(\tau) + \mathbf{X}^T \boldsymbol{\beta} \text{ a}$$

$$\boldsymbol{\beta}_{\mathbf{X}}(\tau) = \begin{pmatrix} \beta_0 + F_{\varepsilon}^{-1}(\tau) \\ \boldsymbol{\beta} \end{pmatrix}.$$

Teda jednotlivé kvantily sa budú líšiť iba v intercepte. Čo by sme čakali, pretože v tomto prípade sa budú marginálne distribúcie líšiť iba v posunutí.

Nech  $h$  je rastúca merateľná funkcia. Potom platí

$$F_{h(Y)|Z}^{-1}(\tau) = \mathbf{X}^T \beta \iff F_{Y|Z}^{-1}(\tau) = h^{-1}(\mathbf{X}^T \beta).$$

Pre nás je hlavne zaujímavá implikácia z ľava do prava. Zaujímavý príklad pre  $h = \log$  potom nech

$$F_{\log(Y)|Z}^{-1}(\tau) = \mathbf{X}^T \beta$$

a máme teda

$$F_{Y|Z}^{-1}(\tau) = \exp(\mathbf{X}^T \beta).$$

Riešime úlohu

$$\min_{\mathbf{b} \in \mathbb{R}^p} \sum_{i=1}^n \rho_{\tau}(Y_i - \mathbf{X}_i^T \mathbf{b}).$$

Je konvexná?

Tento problém vieme previesť do tvaru úlohy lineárneho programovania:

$$\begin{aligned} \min_{\mathbf{b} \in \mathbb{R}^p, \mathbf{r}^+ \geq \mathbf{0}_n, \mathbf{r}^- \geq \mathbf{0}_n} \quad & \tau \sum_{i=1}^n r_i^+ + (1 - \tau) \sum_{i=1}^n r_i^-, \\ \text{s.t.:} \quad & \sum_{j=1}^p X_{ij} b_j + r_i^+ - r_i^- = Y_i, \quad i = 1, \dots, n. \end{aligned}$$

Túto úlohu už vieme ďalej riešiť štandardnými metódami LP.

Za platnosti modelu vieme ukázať, že platí

$$\sqrt{n}(\hat{\beta}_n(\tau) - \beta_{\mathbf{X}}(\tau)) \xrightarrow{D} N_p(\mathbf{0}, \mathbb{V}),$$

kde

$$\mathbb{V} = (\mathbb{E}[\mathbf{X}\mathbf{X}^T f_{Y|Z}(F_{Y|Z}^{-1}(\tau))])^{-1} \tau(1-\tau) \mathbb{E}[\mathbf{X}\mathbf{X}^T (\mathbb{E}[\mathbf{X}\mathbf{X}^T f_{Y|Z}(F_{Y|Z}^{-1}(\tau))])^{-1}.$$

V prípade homoskedastického modelu vieme značne zjednodušiť na tvar

$$\mathbb{V} = \frac{\tau(1-\tau)}{[f_{\varepsilon}(F_{\varepsilon}^{-1}(\tau))]^2} (\mathbb{E}[\mathbf{X}\mathbf{X}^T])^{-1}.$$

# Odhad variančnej matice za predpokladu homoskedasticity

Hlavný problém je v odhade  $s(\tau) = 1/f_\varepsilon(F_\varepsilon^{-1}(\tau))$ . Koenker navrhuje použiť

$$\hat{s}_n = \frac{\hat{F}_{n\hat{\varepsilon}}^{-1}(\tau + h_n) - \hat{F}_{n\hat{\varepsilon}}^{-1}(\tau - h_n)}{2h_n},$$

kde  $h_n$  je ľubovoľná postupnosť konvergujúca k 0 a

$$\hat{F}_{n\hat{\varepsilon}}^{-1}(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{\hat{\varepsilon}_i(\tau) \leq x\} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{Y_i - \mathbf{x}_i^T \hat{\beta}_n(\tau)\}.$$

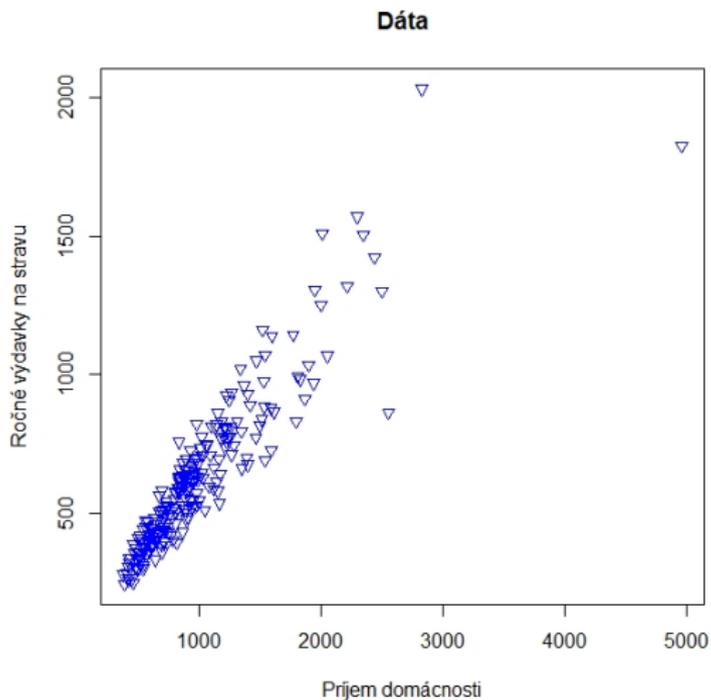
Za predpokladu normálne rozdelených reziduí, Koenker ďalej navrhuje použiť

$$h_n = n^{-1/3} u_{1-\tau/2}^{2/3} \left( \frac{1.5\phi^2(u_\tau)}{2u_\tau^2 + 1} \right)^{1/3},$$

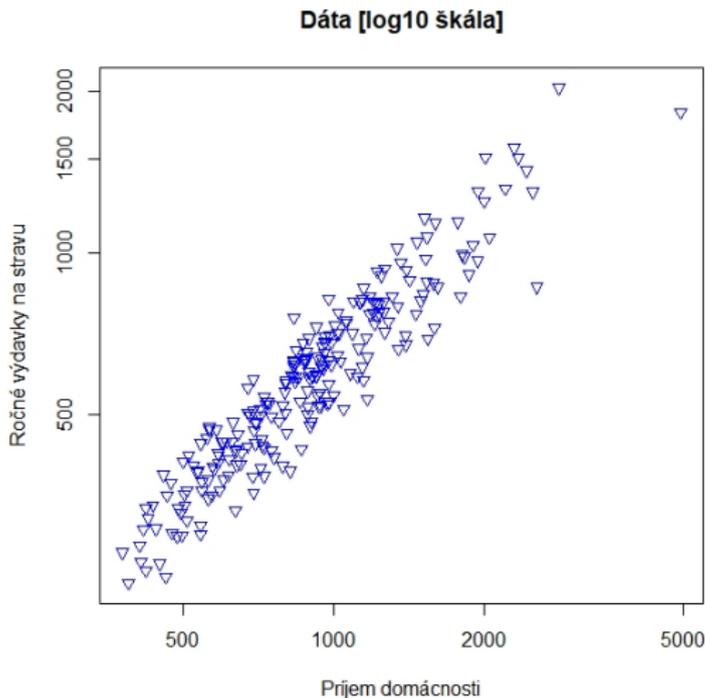
kde  $\phi$  značí hustotu  $N(0, 1)$ .

- Na ilustráciu použijeme dáta "engel" z knižnice "quantreg". Tieto dáta obsahujú ročný príjem (income) a ročné výdavky na jedlo (foodexp) v 235 belgických domácnostiach.
- Táto knižnica taktiež obsahuje funkciu "rq", ktorá je veľmi podobná nám známej funkcii "lm". Na vyrovnanie dát 10%-ným regresným kvantilom použijeme príkaz `'>rq(foodexp~income, tau = 0.1, data = engel)'`.
- "rq" má ako defaultnú hodnotu  $\tau = 0.5$ .
- Jedná z možností je "method", ktorá umožňuje výber výpočetného algoritmu. Pri dátach väčšieho rozsahu sa odporúča zmeniť defaultnú metódu pre urýchlenie výpočtu.

Zrejme homoskedasticita sa nedá predpokladať.

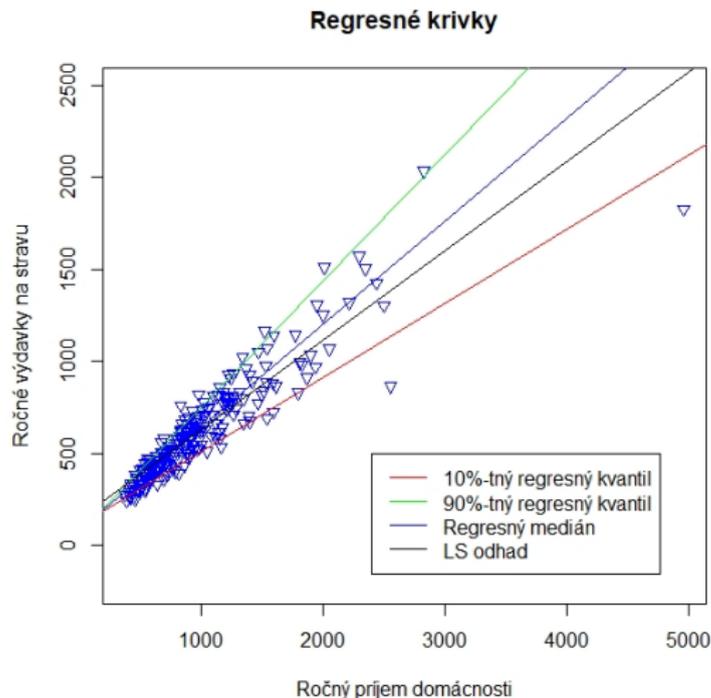


Vieme sa "zbaviť" heteroskedasticity transformáciou odpovede aj vysvetľujúcej premennej.



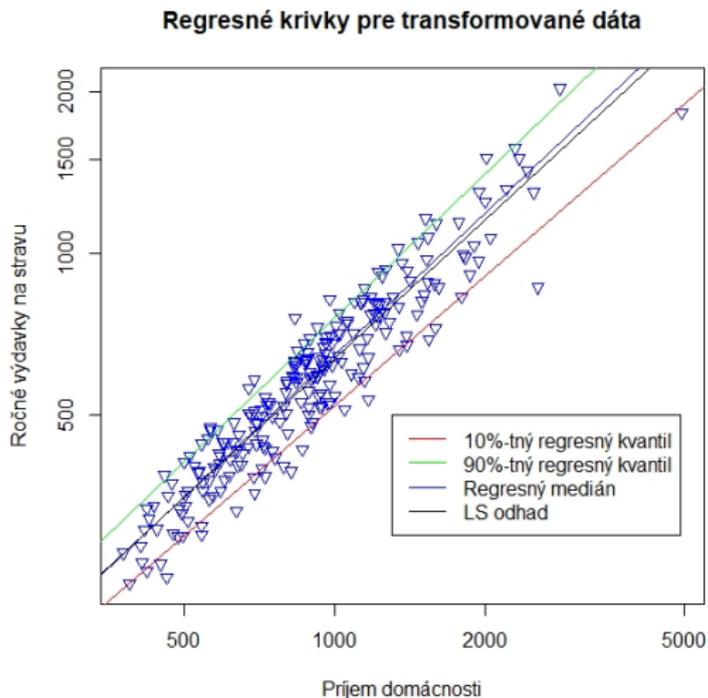
# Odhadnuté hodnoty

- Bez použitia transformácií;
- Priesečník odhadnutých kvantilov - môže byť znamením zlého modelu.



# Odhadnuté hodnoty

- Logaritmická škála;
- Bez priesečníka odhadnutých kvantilov.



- Kvantilová regresia môže byť vhodnejšia z povahy problému;
- Máme radi homoskedastické dáta;
- Tradeoff interpretácia vs. výpočetná zložitosť pri heteroskedastických dátach.

- [1] **Koenker, R.**  
Quantile regression  
*Cambridge University Press, 2005.*
- [2] **Omelka, M.**  
NMST434 - Modern statistical methods, Course notes  
*MFF UK, September, 2021.*