

Univerzita Karlova v Praze  
Matematicko-fyzikální fakulta

## **BAKALÁŘSKÁ PRÁCE**



Peter Rusnák

### **Regresní kvantily**

Katedra pravděpodobnosti a matematické statistiky MFF UK

Vedoucí bakalářské práce: RNDr. Jan Kalina, Ph.D.  
Ústav informatiky AV ČR

Studijní program: matematika  
Studijní obor: finanční matematika

Praha 2011

Ďakujem pánovi RNDr. Janu Kalinovi, Ph.D. za trpezlivé vedenie bakalárskej práce. Ďalej ďakujem svojej rodine, blízkym a každému, kto mi s touto prácou akokoľvek pomohol.

Prohlašuji, že jsem tuto bakalářskou práci vypracoval samostatně a výhradně s použitím citovaných pramenů, literatury a dalších odborných zdrojů.

Beru na vědomí, že se na moji práci vztahují práva a povinnosti vyplývající ze zákona č. 121/2000 Sb., autorského zákona v platném znění, zejména skutečnost, že Univerzita Karlova v Praze má právo na uzavření licenční smlouvy o užití této práce jako školního díla podle § 60 odst. 1 autorského zákona.

V ..... dne.....

podpis

Názov práce: Regresní kvantily

Autor: Peter Rusnák

Katedra / Ústav: Katedra pravděpodobnosti a matematické statistiky MFF UK

Vedúci bakalárskej práce: RNDr. Jan Kalina, Ph.D., Ústav informatiky AV ČR

Abstrakt: Kvantilová regresia je štatistická metóda slúžiaca na určovanie závislostí medzi premennými, ktorá bola navrhnutá už v článku Koenker a Bassett (1978). Od tej doby prešla veľkým rozvojom, keď boli študované jej teoretické vlastnosti, a zároveň si našla radu praktických aplikácií pri spracovaní reálnych dát v najrôznejších oboroch. Kým bežná metóda najmenších štvorcov popisuje vzťah medzi jedným respektíve viacerými kovariátmi  $X$  a podmieneným priemerom odpovedajúcej premennej  $Y$  daným  $X = x$ , kvantilová regresia popisuje vzťah medzi  $X$  a podmienenými kvantilmi  $Y$  danými  $X = x$ . Táto práca obsahuje teóriu nevyhnutnú pre pochopenie vzťahu medzi štandardnou a kvantilovou regresiou a umožňujúcu začlenenie takto získaných odhadov do väčšej skupiny  $M$ -odhadov. Výpočet koeficientu pre jednotlivé kovariáty je prevedený Frisch-Newtnovým algoritmom, ktorý patrí k metódam lineárneho programovania. Taktiež si ukážeme, ako vedľajší produkt tohto algoritmu, takzvané regresné poradie je vypočítané a ako ho použiť pre testovanie hypotéz. V druhej časti budeme ilustrovať numerický výpočet pre kvantilovú regresiu ako na vygenerovaných dátach tak na dátach reálnych.

Kľúčové slová: kvantilová regresia,  $M$ -odhady, regresné poradie, regresné poradové skóry

Title: Regression Quantiles

Author: Peter Rusnák

Department: Department of Probability and Mathematical Statistics

Supervisor: RNDr. Jan Kalina, Ph.D., Institute of Computer Science, AS CR

Abstract: Quantile regression is a statistical method for specifying dependencies among variables, which was introduced by Koenker a Bassett in 1978. Since that time it has gone through a big development, when its theoretical properties have been under study, and it also has found many practical applications for data processing in variety of fields.

While ordinary least-squares regression describes the relationship between one or more covariates  $X$  and the conditional mean of a response variable  $Y$  given  $X = x$ , quantile regression describes the relationship between  $X$  and the conditional quantiles of variable  $Y$  given  $X = x$ . This work contains the theory necessary for understanding relationship between standard and quantile regression and enabling include so received estimates to bigger group of M-estimates. The computation of coefficients for particular covariates is made by using Frisch-Newton algorithm belonging to methods of linear programming. The so-called regression ranks are also obtained as a by-product of this algorithm and we discuss their computational aspects and usage for hypothesis testing. In the second part, we are illustrating numerical computation for quantile regression on theoretical data sets as well as data sets from real world.

Keywords: Quantile regression, M-estimates, regression ranks, regression rank scores

# Obsah

<b>Úvod</b>	<b>1</b>
<b>1. Úvod do kvantilov a regresnej analýzy</b>	<b>2</b>
1.1. Kvantil pravdepodobnostného rozdelenia	2
1.2. Regresná analýza	3
<b>2. M-odhady</b>	<b>5</b>
2.1. Obecný M-odhad	5
2.2. Príklady M-odhadov	6
2.3. M-odhad regresných parametrov	7
2.4. Vlastnosti M-odhadov	7
<b>3. Regresný kvantil</b>	<b>11</b>
3.1. Regresný kvantil ako M-odhad	11
3.2. Regresný kvantil ako riešenie úlohy lineárneho programovania	12
3.3. Regresné poradové skóre	15
<b>4. Aplikácie kvantilovej regresie</b>	<b>19</b>
4.1. Aplikácia prevedená na vygenerovaných dátach	19
4.2. Aplikácia prevedená na praktických dátach	25
<b>Záver</b>	<b>33</b>
<b>Zoznam použitej literatúry</b>	<b>34</b>

## Úvod

Hlavnou úlohou štatistiky je priniesť poriadok do spleti zdanlivo chaotických údajov, k čomu nám slúži analýza dát. Veľmi efektívnym spôsobom, ako analyzovať dáta, je skúmať ich metódami klasickej lineárnej regresie, akou je napríklad najčastejšie používaná metóda najmenších štvorcov. Táto metóda slúži na odhadnutie parametrov lineárneho regresného modelu, prostredníctvom ktorého získame odhad podmienenej strednej hodnoty závislej na hodnotách vysvetľujúcich premenných (regresorov). Stredná hodnota je síce dôležitým ukazovateľom, ktorý predstavuje “stredové” chovanie náhodnej veličiny, avšak nedáva nám žiadnu predstavu o tom, ako sa správa podmienená hustota na krajoch rozdelenia. To priviedlo pánov G. Basseta a R. Koenkera k snahe o zobecnenie obyčajných výberových kvantilov na lineárny regresný model.

Kvantilovú regresiu môžeme použiť na preskúkanie správania sa nezávislých premenných (kovariátov) v každom kvantile skúmanej premennej, čo nám dáva príležitosť pre ucelenejší pohľad na priebeh a vzťahy medzi stochastickými premennými, ako použitie metód klasickej lineárnej regresie.

Táto práca si kladie za úlohu sprostretkovať výklad základných vlastností regresných kvantilov a demonštráciu takto získaných teoretických vedomostí na vhodne zvolených príkladoch. Pred odvodením teórie regresných kvantilov bude nutné zoznámiť sa taktiež so základmi regresie, robustnej štatistiky a optimalizácie.

V prvej kapitole si popíšeme teoretické základy stojace za kvantilmi, regresiou a v krátkosti si teoreticky zhrnieme postup pri odhade metódou najmenších štvorcov. V druhej kapitole si popíšeme základné vlastnosti robustných odhadov, špeciálne sa budeme venovať M-odhadom. V tretej kapitole si zavedieme odhad regresných kvantilov, popíšeme základy lineárneho programovania, ukážeme si ako previesť úlohu odhadnutia regresných koeficientov na úlohu lineárneho programovania a popíšeme si ako ju teoreticky riešiť prostredníctvom Frisch-Newtnovho algoritmu. Taktiež poukážeme na riešenie duálnej úlohy, ktoré nám vznikne pri riešení úlohy lineárneho programovania, toto riešenie nazveme regresné poradové skóre a ukážeme si jeho využitie pri testovaní hypotéz. V poslednej štvrtej kapitole si ukážeme využitie nadobudnutých znalostí na dvoch príkladoch, kde prvý bude demonštratívny na vhodne vygenerovaných dátach a druhý bude predstavovať použitie kvantilovej regresie na reálnych dátach.

# 1. Úvod do kvantilov a regresnej analýzy

## 1.1. Kvantil pravdepodobnostného rozdelenia

Definícia kvantilov vychádza z jednoduchšej optimalizačnej úlohy, ktorá pojednáva o teoretickom probléme bodového odhadu, ktorý je požadovaný pre náhodnú veličinu s distribučnou funkciou  $F$ . Nech strata je popísaná lineárnou funkciou v závislosti od objemu financovaných prostredkov  $x$ :

$$\rho_\tau(x) = x(\tau - I(x < 0)) \text{ pre nejaké } \tau \in (0, 1).$$

Našou úlohou je nájsť  $\tilde{x}$  minimalizujúce očakávanú stratu, a teda budeme sa snažiť minimalizovať

$$E \rho_\tau(X - \tilde{x}) = (\tau - 1) \cdot \int_{-\infty}^{\tilde{x}} (x - \tilde{x}) dF(x) + \tau \cdot \int_{\tilde{x}}^{\infty} (x - \tilde{x}) dF(x).$$

Derivovaním v premennej  $\tilde{x}$  dostávame

$$(1 - \tau) \cdot \int_{-\infty}^{\tilde{x}} dF(x) - \tau \cdot \int_{\tilde{x}}^{\infty} dF(x) = F(\tilde{x}) - \tau.$$

Keďže  $F$  je neklesajúca funkcia, tak každý prvok množiny  $\{x : F(x) = \tau\}$  minimalizuje očakávanú stratu. Ak je táto množina jednoprvková potom  $\tilde{x} = F^{-1}(\tau)$  v opačnom prípade zvolíme za  $\tilde{x}$  najmenší prvok množiny  $\{x : F(x) = \tau\}$ , aby bola dodržaná dohoda, že empirická kvantilová funkcia je zľava spojitá.

**Definícia :** Nech  $F(y) = P(Y \leq y)$  je distribučná funkcia daného rozdelenia pravdepodobnosti náhodnej veličiny  $Y$  a  $\alpha \in (0, 1)$ . Potom kvantilovú funkciu náhodnej veličiny  $Y$  definujeme ako

$$Q(\alpha) = F^{-1}(\alpha) = \inf \{y \in \mathbb{R} : F(y) \geq \alpha\}$$

a číslo  $y_\alpha = Q(\alpha)$  sa nazýva  $\alpha$ -kvantil rozdelenia s distribučnou funkciou  $F(y)$ .

$\alpha$ -kvantil  $y_\alpha$  rozdeľuje definičný obor náhodnej veličiny  $Y$  na dve časti tak, že platí

$$P(Y < y_\alpha) = \alpha \quad \& \quad P(Y \geq y_\alpha) = 1 - \alpha.$$

Kvantily pre niektoré špeciálne hodnoty  $\alpha$  sú označené špeciálnymi názvami a ich hodnoty pre najdôležitejšie rozdelenia sú uvedené v tabuľkách:

- $F^{-1}(\frac{1}{4})$  sa nazýva dolný kvartil a oddeľuje zo štatistického súboru jednu štvrtinu.

- $F^{-1}\left(\frac{1}{2}\right)$  sa nazýva medián a rozdeľuje súbor na dve rovnako početné množiny.
- $F^{-1}\left(\frac{3}{4}\right)$  sa nazýva horný kvantil.

## 1.2. Regresná analýza

Účelom regresnej analýzy je popísať vzťah medzi skúmanou náhodnou veličinou (označme ju  $Y$ ) a známym náhodným vektorom  $\mathbf{Z} = (Z_1, Z_2, \dots, Z_p)$ , ktorý sa na ňu viaže. Naším cieľom bude určiť, na ktorých premenných z  $Z_1, Z_2, \dots, Z_p$  skúmaná premenná  $Y$  závisí, popísať vzťah ako na nich závisí a predpovedať hodnoty  $Y$ , ak máme k dispozícii pozorovania  $\mathbf{Z}$ . Porovnanie  $Y_i$  nemusí priamo závisieť na  $Z_i$ , ale na nejakej jeho transformácii  $\mathbb{X}_i = f(\mathbf{Z}_i)$ ,  $\mathbb{X}_i \in \mathbb{R}^k$ . Náš cieľ dosiahneme prostredníctvom vhodne zvolenej funkcie (označme ju  $g$ ), ktorá je funkciou "vysvetľujúcich premenných"  $X_1, X_2, \dots, X_k$  a vhodne aproximuje skúmanú veličinu, t.j. " $Y = g(\mathbb{X})$ ". V tejto práci sa budeme zaoberať najmä lineárnou regresnou analýzou a teda  $g$  bude mať tvar

$$g(\mathbb{X}) = \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k,$$

kde konštanty  $\beta_1, \dots, \beta_k \in \mathbb{R}$  predstavujú regresné koeficienty.

Model prisluchajúci k takto zvolenej funkcii  $g$  nazveme lineárny regresný model a jeho tvar pre počet pozorovaní  $n > k$  odvodíme ako:

Nech  $x_{11}, \dots, x_{1n}$  ... predstavuje konkrétnych  $n$  hodnôt náhodnej veličiny  $X_1$ .

⋮

Nech  $x_{k1}, \dots, x_{kn}$  ... predstavuje konkrétnych  $n$  hodnôt náhodnej veličiny  $X_k$ .

Nech  $y_1, \dots, y_n$  ... predstavuje konkrétnych  $n$  hodnôt náhodnej veličiny  $Y$ .

Keďže  $g(\mathbb{X})$  je "len vhodnou aproximáciou"  $Y$  tak existujú  $\epsilon_1, \dots, \epsilon_n$  také, že

$$y_1 = \beta_1 x_{11} + \beta_2 x_{21} + \dots + \beta_k x_{k1} + \epsilon_1,$$

$$y_2 = \beta_1 x_{12} + \beta_2 x_{22} + \dots + \beta_k x_{k2} + \epsilon_2,$$

⋮

$$y_n = \beta_1 x_{1n} + \beta_2 x_{2n} + \dots + \beta_k x_{kn} + \epsilon_n.$$

$\epsilon_1, \dots, \epsilon_n$  predstavujú "chyby aproximácie" a pomocou regresnej analýzy sa ich budeme snažiť znižovať prostredníctvom vhodných volení funkcie  $g(\mathbb{X})$ .

Ak o chybovom vektore  $\epsilon = (\epsilon_1, \dots, \epsilon_n)^T$  budeme predpokladať, že jeho zložky sú normálne rozdelené, vzájomne nezávislé náhodné veličiny s nulovou strednou hodno-

tou a rovnakým rozptylom  $\sigma^2$  a teda  $\epsilon \sim N_n(\mathbf{0}_n, \sigma^2 \cdot I_n)$ . Potom hovoríme o klasickom lineárnom regresnom modeli, ktorého tvar môžeme preformulovať ako:

$$E[Y | x_{1_1}, \dots, x_{k_i}] = \beta_1 x_{1_i} + \beta_2 x_{2_i} + \dots + \beta_k x_{k_i}$$

$$\text{var}[Y | x_{1_1}, \dots, x_{k_i}] = \sigma^2$$

Úlohou regresie je získať nejaký odhad  $\hat{\beta}$  vektoru parametrov  $\beta$ , za účelom spočítania  $\hat{y}_i = \mathbf{x}^T \cdot \hat{\beta}$ . Pomocou dobrého odhadu  $\hat{\beta}$  by sme mali získať  $\hat{y}_i$ , čo najbližšie k  $y_i$ . Teraz si ukážeme ako sa odhadujú regresné koeficienty v klasickom lineárnom modeli pomocou najčastejšie používanej metódy najmenších štvorcov.

**Definícia:** Odhad  $\hat{\beta}$  vektoru parametrov  $\beta$  definovaný vzťahom

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^k} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \arg \min_{\beta \in \mathbb{R}^k} (\mathbf{y} - \mathbf{X}\beta)^T \cdot (\mathbf{y} - \mathbf{X}\beta)$$

sa nazýva *odhad metódou najmenších štvorcov*.

Minimalizovaná funkcia sa nazýva reziduálny súčet štvorcov a keďže nie je ohraničená zhora tak  $\hat{\beta}$ , ktoré získame ako riešenie sústavy normálnych rovníc (1), musí byť minimom.

$$\frac{\partial}{\partial \beta} (\mathbf{y} - \mathbf{X}\beta)^T \cdot (\mathbf{y} - \mathbf{X}\beta) = \mathbf{0} \quad (1)$$

Riešenie sústavy bude mať tvar

$$\hat{\beta} = (\mathbf{X}^T \cdot \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

**Veta: Vlastnosti odhadu  $\hat{\beta}$  získaného metódou najmenších štvorcov:**

- 1)  $E \hat{\beta} = \beta$
- 2)  $\text{var} \hat{\beta} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$
- 3)  $\hat{\beta}$  je konzistentný odhad
- 4) Ak  $\epsilon \sim N_n(\mathbf{0}_n, \sigma^2 \cdot I_n)$ , potom  $\sqrt{n} \cdot (\hat{\beta} - \beta) \xrightarrow{d} N_k(\mathbf{0}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$ .

**Dôkaz:** vid' [7] kapitola 13.

**Veta:** Nech platia vyššie uvedené predpoklady a nech navyše platí, že  $Y \sim N_n(\mathbf{X}\beta, \sigma^2 \cdot I_n)$  potom pre reziduálny súčet štvorcov  $SS_e$  platí, že  $\frac{SS_e}{\sigma^2} \sim \chi_{n-k}^2$ .

**Dôkaz:** vid' [7] Veta 13.6.

## 2. M-odhady

Odhady získané lineárnou regresiou sú špeciálnym prípadom *M-odhadov*, ktoré patria k robustným odhadom. Samotná metóda najmenších štvorcov je príliš citlivá k predpokladom normality i k prítomnosti odľahlých pozorovaní, a preto je žiadúce študovať aj iné (robustné) odhady, ktoré by neboli natoľko citlivé. Jednou z možností sú práve *M-odhady*. Zároveň je pravda, že sám odhad metódou najmenších štvorcov patrí do triedy *M-odhadov*. V tejto kapitole uvedieme základy teórie *M-odhadov* umožňujúce začlenenie nami počítaných odhadov do väčšej štatistickej rodiny a ukážeme si dôležité vlastnosti, ktoré na nich môžeme aplikovať.

### 2.1. Obecný M-odhad

Majme štatistický model s náhodným výberom  $y_1, y_2, \dots, y_n$  z neznámej náhodnej veličiny  $Y$ . Označme  $\theta$  ako funkciu distribučnej funkcie náhodnej premennej  $Y$  (štatistický funkcionál)  $\theta = T(P)$ , kde  $P$  je rozdelenie  $Y$ .

*M-odhad* parametru  $\theta$  získame ako riešenie minimalizačnej úlohy

$$\sum_{i=1}^n \rho(y_i, \theta) := \min, \text{ kde } \theta \in \Theta.$$

V tejto bakalárskej práci sa budeme zaoberať príkladom, kde *M-odhad* je model s posunutím  $\theta$  t.j. budeme hľadať riešenie minimalizačnej úlohy

$$\sum_{i=1}^n \rho(y_i - \theta) := \min, \text{ kde } \theta \in \Theta. \quad (2)$$

V prípade ak je funkcia  $\rho(\cdot)$  diferenciovateľná s absolútne spojitou deriváciou, tak môžeme túto minimalizačnú úlohu prepísať do tvaru hľadania riešenia rovnice

$$\sum_{i=1}^n \psi(y_i - \theta) = 0, \text{ kde } \psi(\cdot) = \frac{\partial}{\partial \theta} \rho(\cdot). \quad (3)$$

Predpokladajme, že  $\rho$  je konvexná funkcia, z toho vyplýva, že aj  $\sum_{i=1}^n \rho(y_i - \theta)$  je konvexná funkcia a teda existuje riešenie minimalizačnej úlohy (3). Avšak toto riešenie nemusí byť určené jednoznačne. Napríklad ak  $\rho$  je lineárna funkcia, potom jej derivácia  $\psi$  je konštantná funkcia a teda rovnica (3) nemá riešenie.

Potom riešenie určíme ako

$$\theta := \frac{1}{2} (\theta^+ + \theta^-), \text{ kde } \theta^+ = \inf \left\{ \theta \left| \sum_{i=1}^n \psi(y_i - \theta) < 0 \right. \right\} \quad \theta^- = \sup \left\{ \theta \left| \sum_{i=1}^n \psi(y_i - \theta) > 0 \right. \right\}.$$

## 2.2. Príklady M-odhadov

Z rovností (2) a (3) vyplýva, že *M-odhad* je určený voľbou funkcie  $\rho$ , alebo jej derivácie  $\psi$ , ktorých voľba závisí od vlastností parametra, ktorý chceme odhadnúť. Napríklad, ak má byť parameter zároveň stredom symetrie s okolím v nule, tak je vhodné voliť párnú funkciu  $\rho$  a teda  $\psi$  je tak nepárna funkcia. Tým, že požadujeme aby odhad bol rozbustný, dostávame ohraničenosť influenčnej funkcie a z toho vyplývajúcu ohraničenosť funkcie  $\psi$ . Uvedme si konkrétne príklady takýchto funkcií  $\rho$ .

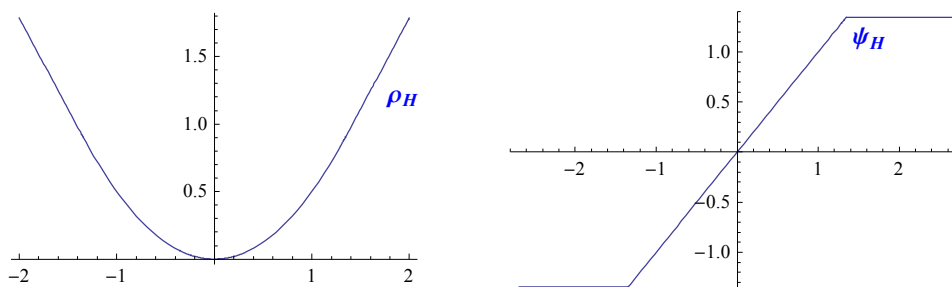
### ■ Huberova funkcia

$$\rho_H(x) = \begin{cases} \frac{x^2}{2} & \dots |x| < k \\ k \cdot \left(|x| - \frac{k}{2}\right) & \dots |x| \geq k \end{cases}$$

kde  $k = c\sigma$  a  $c$  je nejaká konštanta. Empiricky bolo určené, že najvyššiu efektivitu Huberovej funkcie dosiahneme pre  $c=1.345$ .

Príslušná funkcia  $\psi_H$  bude mať tvar:

$$\psi_H(x) = \begin{cases} x & \dots |x| < k \\ k \cdot \text{sgn}(x) & \dots |x| \geq k \end{cases}$$



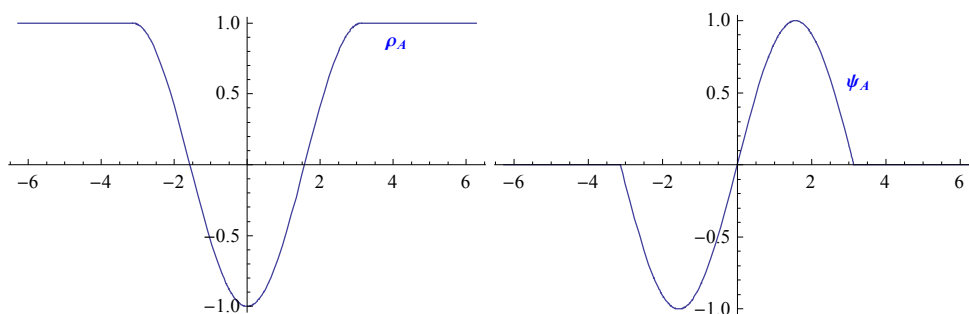
Obrázok 1. Graf Huberovej funkcie a k nej príslušnej funkcie  $\psi_H$

### ■ Andrewsova funkcia

Túto funkciu si uvedieme ako príklad toho, že nie všetky voľby funkcie  $\rho$  sú konvexné, čo spôsobuje, že  $\sum_{i=1}^n \rho(y_i - \theta)$  môže nadobúdať okrem globálnych extrémov aj extrémny lokálne.

$$\rho_A(x) = \begin{cases} -k \cdot \cos \frac{x}{k} & \dots |x| \leq k \cdot \pi \\ k & \dots |x| > k \cdot \pi \end{cases}$$

$$\psi_A(x) = \begin{cases} \sin \frac{x}{k} & \dots |x| \leq k \cdot \pi \\ 0 & \dots |x| > k \cdot \pi \end{cases}$$



Obrázok 2. Graf Andrewsonovej funkcie a k nej príslušnej funkcie  $\psi_A$

### 2.3. M-odhad regresných parametrov

*M-odhad* regresných parametrov je definovaný ako riešenie  $\beta = (\beta_1, \dots, \beta_k)^T$  minimalizujúce súčet funkcií reziduí t.j. rieši minimalizačnú úlohu

$$\sum_{i=1}^n \rho \left( \frac{y_i - \sum_{j=1}^k x_{ij} \beta_j}{\hat{\sigma}} \right) := \min,$$

kde  $\hat{\sigma}$  predstavuje odhad smerodatnej odchýlky.

V prípade ak je funkcia  $\rho(\cdot)$  diferenciovateľná s absolútne spojitou deriváciou, tak môžeme túto minimalizačnú úlohu prepísať do tvaru hľadania riešenia rovnice

$$\sum_{i=1}^n \psi \left( \frac{y_i - \sum_{j=1}^k x_{ij} \beta_j}{\hat{\sigma}} \right) = 0.$$

Existuje veľa používaných funkcií  $\psi$ , ktoré vedú k rôznym M-odhadom napríklad tie uvedené v 2.2.

*M-odhad* s monotónnou funkciou  $\psi$  je doporučený k použitiu, ak náhodné chyby v regresnom modeli majú rozdelenie s chvostami, ktoré sú ťažšie ako chvosty normálneho rozdelenia, ale nie sú ťažšie ako chvosty dvojite exponenciálneho rozdelenia.

### 2.4. Vlastnosti M-odhadov

Napriek tomu, že M-odhady vznikli ako robustné odhady klasických metód, sú robustné len voči odľahlým hodnotám odozvy. Pritom sú veľmi citlivé (nerobustné) voči prítomnosti odľahlých hodnôt v regresoroch (viď [5]), takže má zmysel ich používať len v prípade kontaminácie vysvetľovanej premennej. Vlastnosťami *M-odhadov* budeme rozumieť popísanie vplyvov týchto kontaminácií na hodnotu odhadov.

**Definícia:** Nech  $\mathcal{P}$  je systém všetkých pravdepodobnostných mier na priestore  $(X, \mathcal{B})$ , kde  $X$  je úplný, separabilný metrický priestor a  $\mathcal{B}$  je  $\sigma$ -algebra borelovských podmnožín  $X$  a nech  $P \in \mathcal{P}$ . Povieme, že štatistický funkcionál  $T$  je diferencovateľný v Gateauxovom zmysle podľa  $P$  v smere  $Q$ , ak existuje limita

$$T_Q(P) = \lim_{t \rightarrow 0} \frac{T((1-t)P + t \cdot Q) - T(P)}{t}$$

hodnota tejto limity predstavuje Gateauxovu deriváciu podľa  $P$  v smere  $Q$ .

**Definícia:** Nech  $\epsilon_1, \dots, \epsilon_n$  je náhodný výber z náhodnej veličiny  $Y$ . Potom empirickou (výberovou) distribučnou funkciou pre  $y \in \mathbb{R}$  nazývame funkciu

$$F_n(y) = \frac{1}{n} \sum_{k=1}^n \mathcal{I}_{(-\infty, y)}(y_k).$$

Označme  $P_n$  ako empirickú distribučnú funkciu prisluchajúcu vektoru  $(y_1, \dots, y_n)$ :

$$P_n = \frac{1}{n} \sum_{i=1}^n \delta_{y_i}.$$

Pre ďalšie počítanie v tejto práci budeme predpokladať, že  $P_n$  konverguje k  $P$ , resp. je v istom okolí  $P$  a platí

$$T(P_n) - T(P) = \frac{1}{n} \sum_{i=1}^n T_{\delta_{y_i}}(P) + o_P\left(n^{-\frac{1}{2}}\right) \text{ pre } n \rightarrow \infty.$$

Táto rovnosť nám umožňuje  $\frac{1}{n} \sum_{i=1}^n T_{\delta_{y_i}}(P)$  chápať ako chybu odhadu  $T(P)$  pomocou  $T(P_n)$  a  $T_{\delta_{y_i}}$  ako príspevok  $i$ -teho pozorovania k tejto chybe. To nás privádza k myšlienke definície influenčnej funkcie .

**Definícia:** Gateauxovú deriváciu  $T$  podľa rozdelenia  $P$  v smere  $\delta_y$ ,  $y \in \mathbb{Y}$  nazývame influenčná funkcia funkcionálu  $T$  v rozdelení pravdepodobnosti  $P$  a píšeme

$$IF(y, T, P) = T_{\delta_y}(P) = \lim_{\epsilon \rightarrow 0} \frac{T((1-\epsilon)P + \epsilon\delta_y) - T(P)}{\epsilon}.$$

Pri kvantitatívnych charakteristikách robustnosti nás bude zaujímať, aký veľký vplyv na funkcionál  $T$  má "kontaminácia" rozdelenia  $P$ , ktorá vznikne pridaním nových pozorovaní  $y_{n+1}, y_{n+2}, \dots$  . Pri globálnej citlivosti nás bude zaujímať aký maximálny vplyv by mohlo mať nové pozorovanie  $y_{n+1}$  na funkcionál  $T$ .

**Definícia:** Nech  $T$  je štatistický funkcionál a  $P$  je rozdelenie pravdepodobnosti na  $(\mathcal{Y}, \mathcal{B})$ . Potom globálnou citlivosťou funkcionálu  $T$  pre rozdelenie pravdepodobnosti  $P$  nazývame hodnotu

$$\gamma = \sup_{x \in \mathcal{Y}} \|IF(y, T, P)\|.$$

Pre lepšie pochopenie si uveďme i diskretizovanú formu influenčnej funkcie. Pridajme dodatočné pozorovanie  $y_{n+1}$  k pozorovaniam  $(y_1, \dots, y_n)$ . Potom vplyv tohoto pozorovania môžeme vyjadriť ako

$$I(T_n, y_{n+1}) = T_{n+1}(y_1, y_2, \dots, y_n, y_{n+1}) - T_n(y_1, y_2, \dots, y_n). \quad (4)$$

Citlivosťou funkcionálu  $T_n$  na pridanie ďalšieho pozorovania, tu rozumieme číslo

$$S(T_n) = \sup_y |I(T_n, y)|.$$

Predvedme si to na výpočte globálnej citlivosti strednej hodnoty

$$T(P) = \mathbb{E}_P X, \quad T_n = \bar{X}, \quad T_{n+1} = \bar{X}_{n+1},$$

potom  $T_{n+1}$  môžeme vyjadriť ako

$$T_{n+1} = \frac{1}{n+1} (n \cdot \bar{X} + y)$$

a z (4) tak dostávame

$$I(T_n, y) = \frac{1}{n+1} (n \cdot \bar{X} + y) - \bar{X} = \frac{1}{n+1} (n \cdot \bar{X} - (n+1) \cdot \bar{X} + y) = \frac{1}{n+1} \cdot (y - \bar{X}),$$

prechodom k suprému cez všetky  $Y$  získame

$$S(\bar{X}) = \frac{1}{n+1} \cdot \sup_y |y - \bar{X}| = \infty.$$

A teda globálna citlivosť je rovná nekonečnu, čo znamená, že priemer nie je robustným odhadom strednej hodnoty. Tu sa ukazuje, že M-odhady nemusia mať vždy výhodné robustné vlastnosti. A teda triedu M-odhadov tvoria nielen odhady, ktoré sú robustné, ale patria tam aj iné odhady (napr. priemer), ktoré robustné vlastnosti nemajú. Vo všeobecnosti platí, že ak influenčná funkcia nie je obmedzená, potom odhad príslušného parametru nie je robustný.

Pri lokálnej citlivosti nás bude zaujímať vplyv relatívne malých zmien, pri nahradení pozorovania  $x$  pozorovaním  $y$ , na hodnotu funkcionálu  $T$ .

**Definícia:** Nech  $T$  je štatistický funkcionál a  $P$  je rozdelenie pravdepodobnosti na  $(\mathcal{Y}, \mathcal{B})$ . Potom lokálnou citlivosťou funkcionálu  $T$  pre rozdelenie pravdepodobnosti  $P$  nazývame hodnotu

$$\lambda = \sup_{x,y \in \mathcal{Y}, x \neq y} \frac{\|IF(x,T,P) - IF(y,T,P)\|}{\|x-y\|}.$$

Rovnako ako pri globálnej citlivosti sa pri stanovení odhadu snažíme docieľiť, čo najmenšiu hodnotu lokálnej citlivosti. Avšak pri dosiahnutí veľkých hodnôt (resp. nekonečna) si musíme uvedomiť, že zmena influenčnej funkcie je meraná vzhľadom k rozdielu pozorovaní  $x$  a  $y$  a tak samotná zmena influenčnej funkcie môže byť nepatrná.

Ďalej by nás mohlo zaujímať, koľko minimálne pozorovaní v našom náhodnom výbere by sa muselo zmeniť, aby nový odhad funkcionálu založený na takto upravenom náhodnom výbere bol ľubovoľne vzdialený od nášho pôvodného odhadu.

**Definícia:** Nech  $T$  je štatistický funkcionál a  $P$  je rozdelenie pravdepodobnosti na  $(\mathcal{Y}, \mathcal{B})$ . Označme  $Y^0 = (y_1, \dots, y_n)$  náhodný výber z náhodnej veličiny  $Y$  a k nemu príslušnú hodnotu funkcionálu  $T(Y^0)$ . Nech  $H^m$  je množina všetkých  $m$ -prvkových podmnožín množiny  $\{1, \dots, n\}$ . Označme:

$$\mathcal{Y}_n^m = \{Y_n^m = (\hat{y}_1, \dots, \hat{y}_n) \mid \hat{y}_i = y_i, \text{ pre } i \in \{1, \dots, n\} \setminus h, \hat{y}_i \in \mathcal{Y} \text{ je ľub. pre } i \in h, h \in H^m\}$$

Ďalej označme  $m(T, Y^0)$  najmenšie celé číslo splňujúce

$$\sup_{Y_n^m \in \mathcal{Y}_n^m} \|T(Y_n^m) - T(Y^0)\| = 0.$$

Potom bodom zlyhania funkcionálu  $T$  vo výbere  $Y^0$  nazývame číslo

$$\epsilon_n(T, Y^0) = \left( \frac{m(T, Y^0)}{n} \right). \quad (5)$$

Ak číslo  $m(T, Y^0)$  nezávisí na počiatočnom náhodnom výbere  $Y^0$  potom bodom zlyhania budeme rozumieť limitu

$$\epsilon(T) = \lim_{n \rightarrow \infty} \epsilon_n(T).$$

Bod zlyhania nám tak určuje minimálny podiel potrebných nových pozorovaní. Vo všeobecnosti požadujeme, aby tento pomer bol čo najväčší (ale maximálne  $\frac{1}{2}$ ). Z príkladu o strednej hodnote a z rovnosti (5) dostávame, že ak je influenčná funkcia neobmedzená, tak bod zlyhania je rovný nule.

Existujú aj iné robustné odhady parametrov v lineárnom regresnom modeli ako len *M-odhady*. Ako príklad uveďme *R-odhady* (viď [6]), ktoré sa získajú inverziou testov hypotéz o regresných parametroch založených na poradí (regression ranks). Bližšie sa im budem venovať v kapitole 3.3.

### 3. Regresné kvantily

#### 3.1. Regresný kvantil ako M-odhad

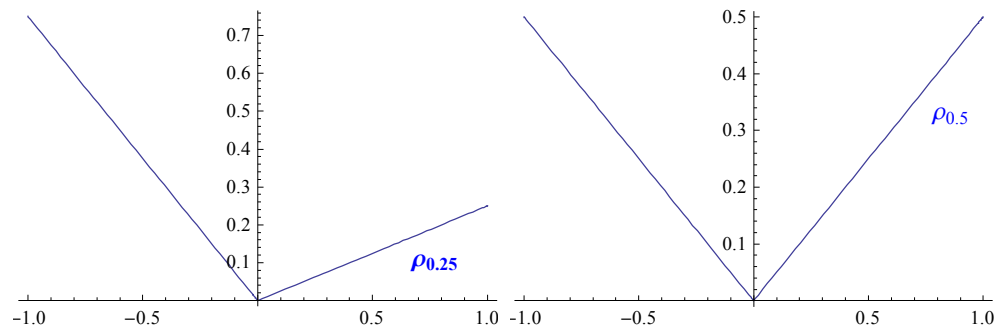
Regresný  $\alpha$ -kvantil pre  $\alpha \in (0,1)$  potom môžeme definovať ako jedno\* z riešení minimalizačnej úlohy, kde za stratovú funkciu  $\rho$  vezmeme tzv. check function:

$$\rho_\alpha(u) = \begin{cases} \alpha \cdot u & \dots u \geq 0 \\ (\alpha - 1) \cdot u & \dots u < 0 \end{cases}$$

respektíve si ju môžeme vyjadriť ako

$$\rho_\alpha(u) = (\alpha - 1) \cdot u \cdot \mathcal{I}_{(-\infty, 0)}(u) + \alpha \cdot u \cdot \mathcal{I}_{[0, \infty)}(u) = u \cdot (\alpha - \mathcal{I}_{(-\infty, 0)}(u)).$$

Pre grafické znázornenie si môžeme vykresliť stratové funkcie pre časté voľby  $\alpha$ :



Obrázok 3. Grafy stratových funkcií  $\rho$  pre  $\alpha=0.25$  resp.  $\alpha=0.5$

Minimalizačná úloha tak bude mať tvar

$$\min_{\beta \in \mathbb{R}^k} \sum_{i=1}^n \rho_\alpha(y_i - x_i^T \beta), \quad (6)$$

čo si môžeme charakterizovať ako súčet reziduí prenasobených  $\alpha$  v prípade, že sú kladné a  $1-\alpha$  ak sú záporné. Táto minimalizačná úloha má aspoň jedno riešenie, pretože platí

$$\lim_{(y_i - x_i^T \beta) \rightarrow \infty} \sum_{i=1}^n \rho_\alpha(y_i - x_i^T \beta) = \lim_{(y_i - x_i^T \beta) \rightarrow -\infty} \sum_{i=1}^n \rho_\alpha(y_i - x_i^T \beta) = \infty.$$

\* minimalizačná úloha môže mať viac riešení ako len jedno, keďže minimalizujeme súčet po častiach lineárnych, konvexných, spojitých funkcií a teda tento súčet je po častiach lineárna, konvexná a spojitá funkcia .

### 3.2. Regresný kvantil ako riešenie úlohy lineárneho programovania

Lineárne programovanie sa zaoberá riešením úloh na viazané extrémny lineárnych funkcií na množinách, ktoré sú určené viazobnými podmienkami v podobe lineárnych rovníc a nerovnic.

Našou úlohou bude nájsť

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^k} \sum_{i=1}^n \rho_{\alpha}(y_i - x_i^T \beta).$$

Majme teda minimalizačnú úlohu

$$\min_{\beta \in \mathbb{R}^k} \sum_{i=1}^n \rho_{\alpha}(y_i - x_i^T \beta).$$

Túto úlohu si môžeme previesť na úlohu lineárneho programovania

$$\min_{(\beta, \epsilon^-, \epsilon^+)} \left\{ \alpha \cdot \mathbf{1}_n^T \cdot \epsilon^+ + (1 - \alpha) \cdot \mathbf{1}_n^T \cdot \epsilon^- \mid \mathbb{X}\beta + \epsilon^+ - \epsilon^- = \mathbf{y}, (\epsilon^+, \epsilon^-) \in \mathbb{R}_+^{2n} \right\}. \quad (7)$$

Zaveďme si označenia:

$\epsilon^-$  pre zápornú časť vektoru residua  $\mathbf{y} - \mathbb{X}\beta$  a  $\epsilon^+$  pre kladnú časť residua  $\mathbf{y} - \mathbb{X}\beta$  t.j.

$$\epsilon_i^- := \begin{cases} -\epsilon_i & \dots \epsilon_i < 0 \\ 0 & \dots \epsilon_i > 0 \end{cases}; \quad \epsilon_i^+ := \begin{cases} \epsilon_i & \dots \epsilon_i > 0 \\ 0 & \dots \epsilon_i < 0 \end{cases} \quad \text{pre } i = 1 \dots n$$

$\beta^-$  pre zápornú časť vektoru  $\beta$  a  $\beta^+$  pre kladnú časť vektoru  $\beta$  t.j.

$$\beta_i^- := \begin{cases} -\beta_i & \dots \beta_i < 0 \\ 0 & \dots \beta_i > 0 \end{cases}; \quad \beta_i^+ := \begin{cases} \beta_i & \dots \beta_i > 0 \\ 0 & \dots \beta_i < 0 \end{cases} \quad \text{pre } i = 1 \dots p$$

Potom môžeme úlohu (6) previesť do tvaru

$$\min_{\mathbf{x}} \{ \mathbf{c}^T \cdot \mathbf{x} \mid \mathbb{A} \cdot \mathbf{x} = \mathbf{y}, \mathbf{x} \geq \mathbf{0} \},$$

kde:

$$\begin{aligned} \mathbf{x} &= (\beta^{+T}, \beta^{-T}, \epsilon^{+T}, \epsilon^{-T}), \\ \mathbf{c} &= (0_p^T, 0_p^T, \alpha \cdot \mathbf{1}_n^T, (1 - \alpha) \cdot \mathbf{1}_n^T), \\ \mathbb{A} &= (\mathbb{X} \quad -\mathbb{X} \quad \mathbf{I}_n \quad -\mathbf{I}_n). \end{aligned}$$

Majme dvojicu úloh

$$\min \{ \mathbf{c}^T \cdot \mathbf{x} \mid \mathbb{A} \cdot \mathbf{x} = \mathbf{y}, \mathbf{x} \geq \mathbf{0} \},$$

$$\max \{ \mathbf{b}^T \cdot \mathbf{y} \mid \mathbb{A} \cdot \mathbf{y} \leq \mathbf{c} \}.$$

Prvú z týchto úloh nazveme primárna, kde lineárna funkcia  $\mathbf{c}^T \cdot \mathbf{x}$  predstavuje účelovú funkciu. Každý vektor  $\mathbf{x}$  splňujúci obmedzenia primárnej úlohy je prípustným riešením minimalizačnej úlohy. Prípustné riešenie  $\mathbf{x}^*$  splňujúce pre každé prípustné riešenie  $\mathbf{x}$ :  $\mathbf{c}^T \cdot \mathbf{x}^* \leq \mathbf{c}^T \cdot \mathbf{x}$  nazveme optimálnym riešením minimalizačnej úlohy. Druhá úloha predstavuje duálnu úlohu k primárnej a platí:

**Veta o dualite:** Majme dvojicu vzájomne duálnych úloh. Potom platí práve jedno z nasledujúcich troch tvrdení:

1. Obe úlohy majú optimálne riešenie a optimálne hodnoty účelových funkcií sú si rovné.
2. Jedna z úloh nemá žiadne prípustné riešenie, druhá úloha má prípustné riešenie, ale nemá optimálne, pretože jej účelová funkcia je na množine prípustných riešení neohraničená.
3. Ani jedna z úloh nemá prípustné riešenie.

**Dôkaz:** vid' [9] Tvrdenie 1.7.

Duálna úloha k našej minimalizačnej úlohe bude mať tvar

$$\max_{\mathbf{b}} \{ \mathbf{b}^T \cdot \mathbf{y} \mid \mathbf{b}^T \cdot \mathbb{X} = \theta, \mathbf{b} \in [\alpha - 1, \alpha]^n \} \quad (\mathbf{b}_i = \alpha \vee \mathbf{b}_i = 1 - \alpha),$$

čo si môžeme prepísať do tvaru :

$$\max_{\mathbf{a}} \{ \mathbf{a}^T \cdot \mathbf{y} \mid \mathbf{a}^T \cdot \mathbb{X} = (1 - \alpha) \mathbf{1}_n^T \cdot \mathbb{X}, \mathbf{a} \in [0, 1]^n \}, \text{ kde } \mathbf{a} = \mathbf{b} + (1 - \alpha) \cdot \mathbf{1}_n^T. \quad (8)$$

Takýto tvar duálnej úlohy zodpovedá tvaru úlohy lineárneho programovania pre použitie metódy vnútorných bodov. Špeciálne vhodnou metódou riešenia je použitie Frisch-Newtnovho algoritmu.

Zaveďme si premennú s splňujúcu :

$$\mathbf{s} = \mathbf{1}_n - \mathbf{a} \quad \text{t.j. } s_i = 1 - a_i \text{ pre } i = 1, \dots, n$$

Teraz preformujeme zadanie našej maximalizačnej úlohy (7), pomocou bariérovej funkcie

$$B(\mathbf{a}, \mathbf{s}, \sigma) = \mathbf{a}^T \cdot \mathbf{y} + \sigma \cdot \sum_{i=1}^n (\ln a_i + \ln s_i).$$

Získame zadanie v tvare

$$\max_{\mathbf{a}} \{ B(\mathbf{a}, \mathbf{s}, \sigma) \mid \mathbf{a}^T \cdot \mathbb{X} = (1 - \alpha) \mathbf{1}_n^T \cdot \mathbb{X}, \mathbf{s} = \mathbf{1}_n - \mathbf{a} \}. \quad (9)$$

Vyriešením tohto problému a následným vypočítaním limity pre  $\sigma \rightarrow 0$  dostaneme riešenie úlohy (7).

**Definícia:** Majme funkciu  $n$  premenných  $f(\mathbf{x}) = f(x_1, \dots, x_n)$ . Gradient  $\nabla f$  definujeme ako vektor parciálnych derivácií funkcie  $f$  vzhľadom k jednotlivým premenným.

Ďalej definujeme Hessián  $\nabla^2 f$  ako maticu druhých parciálnych derivácií.

$$\nabla f = \left( \frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)$$

$$\nabla^2 f = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_n} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

Vyjadrením gradientu a Hessiánu Bariérovej funkcie dostaneme:

$$\frac{\partial B}{\partial a_i} = y_i + \sigma \left( \frac{1}{a_i} - \left( \frac{1}{1-a_i} \right) \right) \quad \Rightarrow \nabla B(\mathbf{a}, \mathbf{s}, \sigma) = \mathbf{y} + \sigma(\mathbf{A}^{-1} - \mathbf{S}^{-1}) \cdot \mathbf{1}_n$$

$$\text{kde } \mathbf{A} = \text{diag}(a_1, \dots, a_n); \mathbf{S} = \text{diag}(s_1, \dots, s_n)$$

$$\frac{\partial^2 B}{\partial a_i \partial a_j} = \begin{cases} \sigma \left( -\frac{1}{a_i^2} - \left( \frac{1}{1-a_i} \right)^2 \right) & \text{pre } i = j \\ 0 & \text{pre } i \neq j \end{cases} \Rightarrow \nabla^2 B(\mathbf{a}, \mathbf{s}, \sigma) = -\sigma(\mathbf{A}^{-2} - \mathbf{S}^{-2})$$

V každom iteračnom kroku budeme chcieť nájsť také  $\mathbf{p} = \mathbf{a}^{(k+1)} - \mathbf{a}^{(k)}$ , ktoré maximalizuje hodnotu

$$B(\mathbf{a}^{(k)} + \mathbf{p}, \mathbf{s}^{(k)} - \mathbf{p}, \sigma).$$

Pre zjednodušenie budeme odteraz označovať  $\mathbf{a}^{(k)}$  len ako  $\mathbf{a}$ . Vyjadríme Taylorov rozvoj druhého rádu barierovej funkcie v premennej  $\mathbf{p}$  ako

$$B(\mathbf{a} + \mathbf{p}, \mathbf{s} - \mathbf{p}, \sigma) = B(\mathbf{a}, \mathbf{s}, \sigma) + \mathbf{p}^T \cdot \nabla B(\mathbf{a}, \mathbf{s}, \sigma) + \frac{1}{2} \cdot \mathbf{p}^T \cdot \nabla^2 B(\mathbf{a}, \mathbf{s}, \sigma) \cdot \mathbf{p}.$$

Podmienku z (9) v tvare

$$(\mathbf{a} + \mathbf{p})^T \cdot \mathbf{X} = \mathbf{a}^T \cdot \mathbf{X} + \mathbf{p}^T \cdot \mathbf{X} = (1 - \alpha) \mathbf{I}_n^T \cdot \mathbf{X},$$

prevedieme do tvaru

$$\mathbf{p}^T \cdot \mathbf{X} = 0.$$

Vďaka predpokladu, že  $\mathbf{a}$  je prípustné riešenie, tak má platiť

$$\mathbf{a}^T \cdot \mathbf{X} = (1 - \alpha) \mathbf{I}_n^T \cdot \mathbf{X}.$$

Z Taylorovho rozvoja, prevedenej podmienky a úlohy (9) tak dostaneme úlohu v tvare

$$\max_{\mathbf{p}} \left\{ B(\mathbf{a}, \mathbf{s}, \sigma) + \mathbf{p}^T \cdot \mathbf{y} + \mathbf{p}^T \cdot \sigma(\mathbf{A}^{-1} - \mathbf{S}^{-1}) \cdot \mathbf{1}_n - \frac{1}{2} \cdot \mathbf{p}^T \cdot \sigma(\mathbf{A}^{-2} - \mathbf{S}^{-2}) \cdot \mathbf{p} \mid \mathbf{p}^T \cdot \mathbf{X} = 0 \right\}.$$

Keďže  $B(\mathbf{a}, \mathbf{s}, \sigma)$  nezávisí na  $\mathbf{p}$  dostávame tvar (nás totiž zaujíma  $\text{argmax}_{\mathbf{p}}$ )

$$\max_{\mathbf{p}} \left\{ \mathbf{p}^T \cdot \mathbf{y} + \mathbf{p}^T \cdot \sigma(\mathbf{A}^{-1} - \mathbf{S}^{-1}) \cdot \mathbf{1}_n - \frac{1}{2} \cdot \mathbf{p}^T \cdot \sigma(\mathbf{A}^{-2} - \mathbf{S}^{-2}) \cdot \mathbf{p} \mid \mathbf{p}^T \cdot \mathbf{X} = 0 \right\}.$$

Vyriešením takéhoto problému a posunutím z  $\mathbf{a}^{(k)}$  (naše  $\mathbf{a}$ ) vo výslednom smere  $\mathbf{p}$ , predstavujúcim smer k hranici množiny prípustných riešení, obdržíme opäť prípustné riešenie a hodnota účelovej funkcie sa nezmenší. Túto úlohu budeme ďalej riešiť prostredníctvom Lagrangerovej metódy neurčitých koeficientov. Nech  $\mathbf{d}$  predstavuje vektor Lagrangerových multiplikátorov, potom Lagrangerova funkcia bude mať tvar

$$L(\mathbf{p}, \mathbf{d}) = \mathbf{p}^T \cdot \mathbf{y} + \mathbf{p}^T \cdot \sigma(\mathbf{A}^{-1} - \mathbf{S}^{-1}) \cdot \mathbf{1}_n - \frac{1}{2} \cdot \mathbf{p}^T \cdot \sigma(\mathbf{A}^{-2} - \mathbf{S}^{-2}) \cdot \mathbf{p} - \mathbf{p}^T \cdot \mathbf{X} \cdot \mathbf{d}.$$

Položíme parciálne derivácie podľa  $\mathbf{p}$  a  $\mathbf{d}$  rovné nule a dostávame:

$$\mathbf{y} + \sigma(\mathbf{A}^{-1} - \mathbf{S}^{-1}) \cdot \mathbf{1}_n - \sigma(\mathbf{A}^{-2} - \mathbf{S}^{-2}) \cdot \mathbf{p} - \mathbf{X} \cdot \mathbf{d} = 0$$

$$\mathbf{p}^T \cdot \mathbf{X} = 0$$

Predpokladáme, že  $\sigma \rightarrow 0$  a teda z prvej rovnice dostávame  $\mathbf{y} \rightarrow \mathbf{X} \cdot \mathbf{d}$  z čoho môžeme predpokladať, že  $\mathbf{d} \rightarrow \beta$ .

Položme  $\mathbf{W} = (\mathbf{A}^{-2} - \mathbf{S}^{-2})^{-1}$ , potom pre násobením prvej rovnice výrazom  $\mathbf{X}^T \cdot \mathbf{W}$  a použitím druhej rovnosti za účelom eliminácie  $\mathbf{p}$  dostaneme

$$\mathbf{X}^T \cdot \mathbf{W} \cdot \mathbf{y} + \sigma \cdot \mathbf{X}^T \cdot \mathbf{W} (\mathbf{A}^{-1} - \mathbf{S}^{-1}) \cdot \mathbf{1}_n = \mathbf{X}^T \cdot \mathbf{W} \cdot \mathbf{X} \cdot \mathbf{d}.$$

Odkiaľ môžeme explicitne vyjadriť vektor Lagrangerových multiplikátorov  $\mathbf{d}$  ako

$$\mathbf{d} = (\mathbf{X}^T \cdot \mathbf{W} \cdot \mathbf{X})^{-1} \cdot [\mathbf{X}^T \cdot \mathbf{W} \cdot \mathbf{y} + \sigma \cdot \mathbf{X}^T \cdot \mathbf{W} (\mathbf{A}^{-1} - \mathbf{S}^{-1}) \cdot \mathbf{1}_n]$$

a teda:

$$\hat{\beta} = \lim_{\sigma \rightarrow 0} \left( (\mathbf{X}^T \cdot \mathbf{W} \cdot \mathbf{X})^{-1} \cdot [\mathbf{X}^T \cdot \mathbf{W} \cdot \mathbf{y} + \sigma \cdot \mathbf{X}^T \cdot \mathbf{W} (\mathbf{A}^{-1} - \mathbf{S}^{-1}) \cdot \mathbf{1}_n] \right).$$

Čo už sme schopný aritmeticky dopočítať a získať tak  $\hat{\beta}$  odhad parametra  $\beta$  lineárneho modelu  $Y = \beta^T \cdot \mathbf{X}$  pre  $\alpha$ -regresný kvantil.

### 3.3. Regresné poradové skóre

Koenker a Basset (1978) charakterizovali regresné kvantily ako riešenie parametrickej úlohy lineárneho programovania. Zodpovedajúca množina duálnych úloh pre nich predstavovala len technický nástroj využitý pri výpočte regresných kvantilov. Teraz ukážeme, že duálne riešenia, ktoré budeme nazývať regresné poradia (regression ranks) majú štatistický význam a môžu byť, veľmi prirodzeným spôsobom, použité na definovanie regresných poradových skór (regression rank scores) v lineárnych modeloch, umožňujúc nám tak mnoho aplikácii medzi, ktoré patrí vytvorenie intervalu spoľahlivosti pre  $\beta_i$  a testovanie hypotéz pre  $\beta_i$  (viď [2]).

Majme duálnu úlohu

$$\operatorname{argmax}_{\mathbf{a}} \{ \mathbf{a}^T \cdot \mathbf{y} \mid \mathbf{a}^T \cdot \mathbf{X} = (1 - \alpha) \mathbf{1}_n^T \cdot \mathbf{X}, \mathbf{a} \in [0, 1]^n \}.$$

Našou úlohou bude skúmať štatistické vlastnosti riešenia  $\hat{\mathbf{a}}(\alpha) = (\hat{a}_1(\alpha), \dots, \hat{a}_n(\alpha))$  tejto úlohy, ktoré si teraz charakterizujeme prostredníctvom teórie lineárneho programovania. Nech  $M = \{i \in \{1, \dots, n\} \mid y_i = \mathbf{x}_i^T \beta(\alpha)\}$ , potom môžeme  $\mathbf{a}(\alpha)$  určiť prostredníctvom nerovnic:

$$\hat{\mathbf{a}}_i(\alpha) = \begin{cases} 1 & \text{ak } y_i > \mathbf{x}_i^T \beta(\alpha) \\ 0 & \text{ak } y_i < \mathbf{x}_i^T \beta(\alpha) \end{cases} \text{ pre } i = 1, \dots, n,$$

a pomocou lineárnych rovníc:

$$\sum_{i \in M} \hat{\mathbf{a}}_i(\alpha) \cdot \mathbf{x}_i = (1 - \alpha) \cdot \sum_{i=1}^n \mathbf{x}_i - \sum_{i=1}^n I[Y_i > \mathbf{x}_i^T \beta(\alpha)] \cdot \mathbf{x}_i.$$

Pripomeňme si, že takéto riešenie existuje na základe *Vety o dualite* vždy, ak existuje riešenie úlohy (6). Ak má navyše "chyba aproximácie"  $\epsilon$  spojitú distribučnú funkciu, tak toto riešenie je jednoznačné pre každé  $\alpha \in (0,1)$ . Výpočtu sme sa venovali v predchádzajúcej kapitole.

Teraz si zavedieme tzv. skórovú funkciu (score function)  $\varphi: (0,1) \rightarrow \mathbb{R}$ , o ktorej budeme predpokladať, že má konečnú variáciu a je konštantná mimo intervalu  $[\epsilon, 1-\epsilon]$  pre nejaké  $\epsilon \in (0,1)$ . Pomocou tejto skórovej funkcie budeme definovať skóry  $\hat{b}_i$  (regression rank scores) ako

$$\hat{b}_i = - \int_0^1 \varphi(s) d\hat{\mathbf{a}}_i(s) \quad i = 1, \dots, n.$$

kde  $\hat{\mathbf{a}}_i(s)$ ,  $s \in (0,1)$  je regresné poradie.

Ako príklad si uvedieme klasickú Wilcoxonovu skórovú funkciu centrovanú v 0:

$$\varphi(\alpha) = \alpha - \frac{1}{2},$$

ktorej integrovanie po častiach nám dáva:

$$\begin{aligned} \hat{b}_i &= - \int_0^1 \left(s - \frac{1}{2}\right) d\hat{\mathbf{a}}_i(s) \\ &= \int_0^1 \hat{\mathbf{a}}_i(s) s - \frac{1}{2} \\ &= \sum_{j=1}^J \frac{1}{2} (\hat{\mathbf{a}}_i(\alpha_{j+1}) - \hat{\mathbf{a}}_i(\alpha_j)) - \frac{1}{2}, \end{aligned}$$

kde  $\alpha_1, \dots, \alpha_J$  sú body zlomu funkcie  $\hat{\mathbf{a}}_i(\alpha)$  a  $\alpha_{J+1} = 1$ .

Pre takúto a podobné voľby skórovej funkcie  $\varphi$  môže byť skonštruovaná testová štatistika  $T_n$ , ktorú vieme použiť na skúmanie globálneho efektu regresorov v matici  $\mathbb{Z}$  na pozorovanú veličinu  $Y$ .

Rozdelíme regresory do dvoch skupín na regresory v matici  $\mathbb{X}$  prostredníctvom, ktorých budeme odhadovať pozorovanú veličinu  $Y$  a regresory v matici  $\mathbb{Z}$ , u ktorých budeme skúmať či ich prídanie do regresného modelu malo efekt na odhad pozorovanej veličiny.

Budeme uvažovať model

$$Y = \mathbb{X} \cdot \beta(\alpha) + \mathbb{Z} \cdot \zeta(\alpha) + \epsilon.$$

Splňujúci:

$$Y_i = \beta_1(\alpha) + x_{i_2} \beta_2(\alpha) + \dots + x_{i_{k-q}} \beta_{k-q}(\alpha) + z_{i_1} \zeta_1(\alpha) + \dots + z_{i_q} \zeta_q(\alpha) + \epsilon_i,$$

$\mathbf{Z} \in \mathbb{R}^{n \times q}$  je matica, ktorej riadky tvoria trojuholníkové pole,

na ktorom budeme testovať nulovú hypotézu  $H_0: \zeta(\alpha) = o$ .

Spočítame regresné poradové skóry z "restricted form" (použijeme len maticu  $\mathbb{X}$ ):

$$\operatorname{argmax}_{\mathbf{a}} \{ \mathbf{a}^T \cdot \mathbf{y} \mid \mathbf{a}^T \cdot \mathbb{X} = (1 - \alpha) \mathbf{1}_n^T \cdot \mathbb{X}, \mathbf{a} \in [0, 1]^n \}.$$

Skúmaná testová štatistika má tvar

$$T_n = S_n^T M_n^{-1} S_n / A^2(\varphi), \quad (10)$$

kde:

$$S_n = n^{-1/2} (\mathbf{Z} - \hat{\mathbf{Z}}) \hat{\mathbf{b}},$$

$$\mathbf{Z} - \hat{\mathbf{Z}} = \mathbf{Z} - \mathbb{X} (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbf{Z},$$

$$M_n = (\mathbf{Z} - \hat{\mathbf{Z}}) (\mathbf{Z} - \hat{\mathbf{Z}})^T / n,$$

$$\hat{\mathbf{b}} = \int_0^1 \hat{\mathbf{a}}(s) d\varphi(s) = (\hat{b}_1, \dots, \hat{b}_n),$$

$$A^2(\varphi) = \int_0^1 (\varphi(s)^2 - \int_0^1 \varphi(s) ds)^2 ds.$$

$\mathbf{Z} - \hat{\mathbf{Z}}$  predstavuje reziduá vzniklé odhadom  $\mathbf{Z}$  prostredníctvom metódy najmenších štvorcov vykonanej na dátach  $\mathbb{X}$ .

Pri platnosti vyššie uvedených predpokladov a nulovej hypotézy má  $T_n$  limitné rozdelenie  $\chi_q^2$ , kde  $q$  predstavuje počet zložiek parametra  $\zeta$ . Všimnime si, že takto navrhnutý test za platnosti nulovej hypotézy nie je vôbec závislý na distribučnej funkcii  $F$  "chyby aproximácie"  $\epsilon$  a ani na hodnotách  $\beta$ .

Dôležitou aplikáciou teórie testovania založenej na regresných poradových skóroch je vytváranie intervalov spoľahlivosti pre jednotlivé parametre kvantilového regresného modelu s konkrétnym kvantilom  $\alpha$ . Štatistiku  $T_n$  budeme vytvárať pomocou skórovej funkcie  $\varphi$ , ktorá sa bude exkluzívne sústrediť na jednu hodnotu  $\alpha$  (chceme otestovať vplyv regresorov na hodnotu pozorovaní  $Y$  len v časti populácie zodpovedajúcej  $\alpha$ ).

Budeme uvažovať model

$$Q_{Y_i}(\alpha \mid x_i, z_i) = x_i^T \beta(\alpha) + z_i^T \zeta(\alpha).$$

Naším cieľom bude určiť, či vektor regresorov  $z_i$  má vplyv na sledovanú náhodnú veličinu  $Y$  pre zvolený  $\alpha$ -regresný kvantil.

Testujeme hypotézu  $H_0 : \zeta = \zeta_0$  (špeciálny prípad, ktorý nás najčastejšie bude zaujímať je  $\zeta=0$ ), tento test bude založený na regresných poradiach, ktoré môžeme získať ako riešenie úlohy

$$\operatorname{argmax}_{\mathbf{a}} \{ \mathbf{a}^T \cdot (\mathbf{y} - z\zeta_0) \mid \mathbf{a}^T \cdot \mathbb{X} = (1 - \alpha) \mathbf{1}_n^T \cdot \mathbb{X}, \mathbf{a} \in [0, 1]^n \}.$$

Skórová funkcia

$$\varphi_\alpha(s) = \alpha - I(s < \alpha) \quad (11)$$

nám umožňuje zamerať sa počas testu výhradne na  $\alpha$ -kvantil.

Špeciálny prípad, pre medián majúci tvar

$$\varphi_{1/2}(t) = \frac{1}{2} \operatorname{sgn}\left(t - \frac{1}{2}\right),$$

je dôležitou aplikáciou, ktorá sa často označuje ako sgn-skóre.

Použitím skórovej funkcie (11) vypočítame skóre ako

$$\hat{b}_i = - \int_0^1 \varphi(s) d\hat{\mathbf{a}}_i(s) = \hat{\mathbf{a}}_i(\alpha) - (1 - \alpha) \quad i = 1, \dots, n.$$

Pri platnosti  $H_0$  platí

$$S_n(\zeta_0) = n^{-1/2} (Z - \hat{Z})^T \hat{\mathbf{b}} \xrightarrow{\mathcal{P}} S \sim N(0, A^2(\varphi_\alpha) \cdot q_n^2),$$

kde:

$$q_n^2 = n^{-1} Z^T (I - X(X^T X)^{-1} X^T) Z,$$

$$\hat{\mathbf{b}} = (\hat{b}_1, \dots, \hat{b}_n),$$

$$A^2(\varphi_\alpha) = \int_0^1 (\varphi_\alpha(s) - \int_0^1 \varphi_\alpha(s) ds)^2 ds = \alpha \cdot (1 - \alpha).$$

Potom platí, že

$$T_n(\zeta_0) = S_n(\zeta_0) / (A(\varphi_\alpha) q_n) \xrightarrow{\mathcal{P}} T \sim N(0, A^2(\varphi_\alpha) \cdot q_n^2).$$

Hypotézu  $H_0$  zamietame, ak  $|T_n(\zeta_0)| > \Phi^{-1}(1 - \frac{\alpha}{2})$ .

## 4. Aplikácie kvantilovej regresie

### 4.1. Aplikácia prevedená na vygenerovaných dátach

Pre lepšie pochopenie, toho ako funguje kvantilová regresia, si ukážeme použitie štatistických nástrojov z nej vychádzajúcich na jednoduchých štatistických dátach, ktoré si sami vygenerujeme takým spôsobom, aby výsledky mohli byť jednoducho graficky prezentované. Výpočty prevedieme v štatistickom programe R (viď [10]), niektoré pomocné výpočty si predpripravíme v matematickom programe *Mathematica*.

Zvolíme si regresný model s dvoma regresormi:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 z_i + \epsilon_i, \text{ pre } i = 1, \dots, n. \quad (12)$$

Teraz si vhodne zvolíme hodnoty skutočných regresných koeficientov a to takým spôsobom, aby sme si na nich mohli prezentovať rôzne výsledky pri testovaní hypotéz významnosti jednotlivých regresorov (viď predchádzajúcu kapitolu).

$$(\beta_0, \beta_1, \beta_2) = (1, 2, 0.01)$$

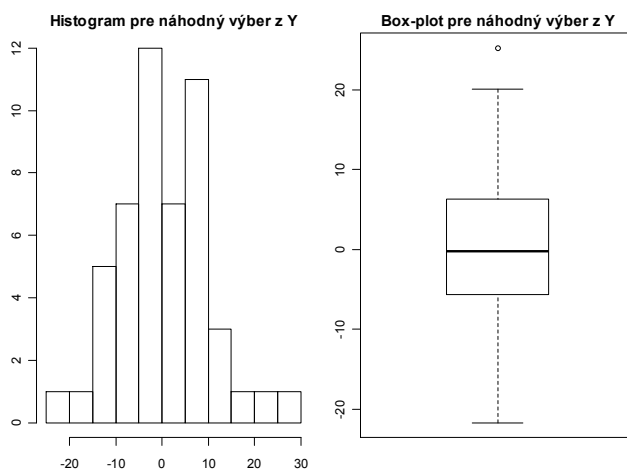
Počet pozorovaní  $n$  si zvolíme 50 a následne si vhodne vygenerujeme hodnoty regresorov a chýb aproximácie:

- hodnoty regresora  $x_i$  budú predstavovať náhodný výber z normálneho rozdelenia  $\mathcal{N}(0,5)$
- hodnoty regresora  $z_i$  budú predstavovať náhodný výber z rovnomerného rozdelenia  $\mathcal{R}(-1,1)$
- hodnoty chýb aproximácie  $\epsilon_i$  predstavujú náhodný výber z normálneho rozdelenia  $\mathcal{N}(0,1)$

Teraz si prostredníctvom takto vygenerovaných hodnôt a ich vložením do modelu (12) vypočítame príslušné hodnoty  $y_1, \dots, y_n$  skúmanej veličiny  $Y$ .

Bližší náhľad na získané hodnoty  $y_i$ :

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
-21.67	-5.5250	-0.2492	0.3470	6.2730	25.2100



Obrázok 4. Histogram a boxplot vygenerovaných hodnôt  $y_i$

Ešte pred samotnou regresnou analýzou si zvolíme, čo najvhodnejší model, na ktorý ju budeme aplikovať. Na začiatku uvažujeme náš pôvodný regresný model (12)

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 z_i + \epsilon_i, \text{ pre } i = 1, \dots, 50.$$

Následne sa budeme snažiť určiť na ktorých regresoroch je skúmaná veličina závislá a naopak, ktoré regresory môžeme za istých okolností z modelu vylúčiť. To môžeme docieľiť viacerými metódami:

- **Procedúra stepwise**

Táto metóda minimalizuje súčet reziduálnych štvorcov prostredníctvom pridávania, alebo odoberania premenných spôsobom, ktorý zohľadňuje štatistickú významnosť všetkých premenných v modeli. Táto stepwise procedúra bude kombinovať prvky napredovacieho (forward) výberu a spätnej (backward) eliminácie, keď po každom eliminačnom štádiu bude nasledovať kontrola pre možné pridanie. Ako počiatočný štatistický objekt reprezentujúci model, ku ktorému procedúra začne pridávať regresory, zvolíme objekt vzniklý kvantilovou regresnou analýzou, ktorá skúmanú veličinu modeluje len pomocou konštanty. Výpočet regresných koeficientov prebieha prostredníctvom Frisch-Newtonovho algoritmu, ktorého priebeh sme si ukázali v predchádzajúcej kapitole. Procedúre umožníme pridávať regresory obsiahnuté v našom modeli (12).

```
model0<-y~x+z
model<-step(rq(y~1,method="fn"),scope=model0,trace=0)$formula
```

procedúra nám vráti:

```
Call:
rq(formula = y ~ x, method = "fn")
```

```
Coefficients:
(Intercept)          x
  1.108684      2.025940
```

```
Degrees of freedom: 50 total; 48 residual
```

Poznamenajme, že program R vo výstupnom modeli neuvádza konštantu, ktorá je však v modeli obsiahnutá. Procedúra navrhuje vynechať z pôvodného modelu regresor Z. Model navrhnutý stepwise procedúrou, ako vhodný na regresnú analýzu, tak má tvar

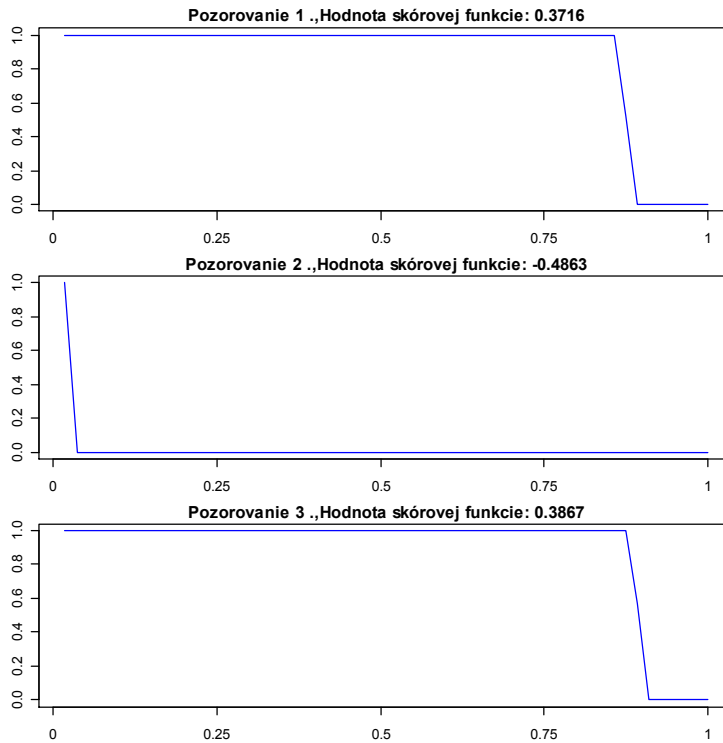
$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ pre } i = 1, \dots, 50. \quad (13)$$

- **Testovanie hypotéz významnosti**

Tento test bude založený na regresných poradniach a pre príslušné “restricted form” nášho modelu ( $Y \sim X$  resp.  $Y \sim Z$ ). Pre lepšie pochopenie toho, čo vlastne sú regresné poradia si ich môžeme v programe R vyjadriť pre “restricted form”  $Y \sim X$  použitím príkazu:

```
model<-y~x
regres2<-rq(model,tau=-1,method="br")$sol
```

a následne graficky znázorniť :



**Obrázok 5.** Grafy regresných poradových skór prvých troch pozorovaní pre  $y \sim x$

Regresné poradové skóry  $\hat{b}$  z príslušnej “restricted form”  $Y \sim Z$ , využité pri výpočte významnosti regresora  $X$ , vypočítame v programe R použitím príkazu

```
skore1<-ranks(rq(y~z,tau=-1,method="br"),score="wilcoxon").
```

Analogicky vypočítame regresné skóry pre “restricted form”  $Y \sim X$ :

```
skore2<-ranks(rq(y~x,tau=-1,method="br"),score="wilcoxon")
```

Ako skórovú funkciu sme zvolili Wilcoxonovu funkciu, ktorá má tvar

$$\varphi(t) = 2t - 1 \quad ; \quad 0 < t < 1.$$

Teraz použijeme testovú štatistiku (10) založenú na regresných poradových skóрах

$$T_n = S_n^T M_n^{-1} S_n / A^2(\varphi), \text{ pre } n = 1, 2,$$

o ktorej na základe teórie uvedenej v kapitole 3.3 predpokladáme, že má  $\chi_1^2$  rozdelenie.

Hodnoty testovej štatistiky a ich *p*-hodnoty si uvedieme v tabuľke:

Regresor	Testová štatistika $T_n$	<i>p</i> -hodnota
$X$	25.1642	5.265054 e – 07
$Z$	0.00372309	0.9513456

**Tabuľka 1.** Testy závislosti na jednom regresore pre TP

Z tabuľky vidíme, že na 5% hladine významnosti nezamietame nevýznamnosť regresora  $Z$ . Rovnako ako stepwise procedúra teda aj testovanie hypotéz významnosti navrhuje vynechať z pôvodného modelu regresor  $Z$ .

Ako optimálny model pre regresnú analýzu nám tak v oboch prípadoch vyšiel model v tvare

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ pre } i = 1, \dots, 50.$$

Na tento model (13) teda budeme aplikovať kvantilovú regresiu prevedenú v kvantiloch použitím príkazu

```
regresory <- summary(rq(model, tau = c(0.25, 0.5, 0.75)))
```

Výsledky odhadov regresných koeficientov spolu s 95 % intervalmi spoľahlivosti si uvedieme v tabuľke:

Regr. koeficient	0.25	0.5	0.75
$\beta_0$	0.25781 (-0.01760, 0.58221)	1.10868 (0.66501, 1.41321)	2.24806 (1.58247, 2.38872)
$\beta_1$	1.98903 (1.89188, 2.07814)	2.02594 (1.90321, 2.08260)	1.96722 (1.92740, 2.09957)

**Tabuľka 2.** Výsledky odhadov regresných koeficientov kvantilovou regresiou pre TP

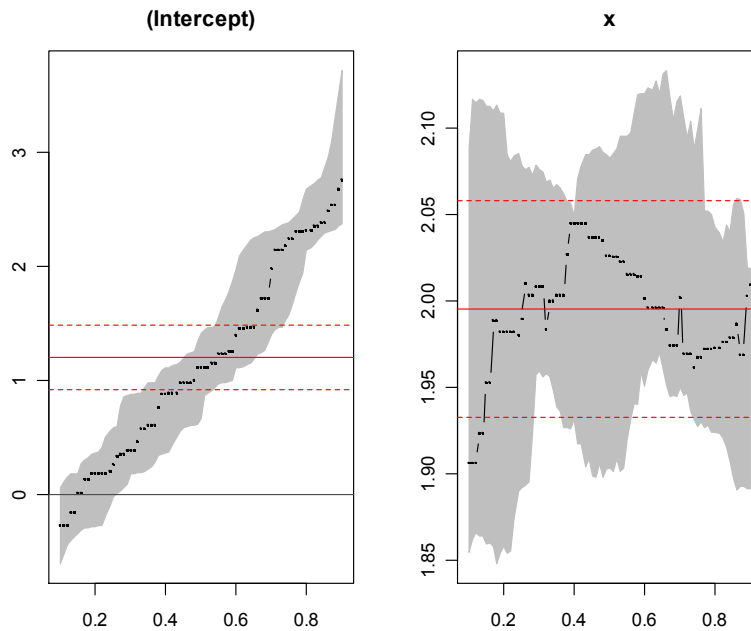
Poznamenajme, že regresný koeficient  $\beta_2$  volíme na základe testov o významnosti regresorov nulový vo všetkých kvantiloch.

Pre úplnosť si ešte uvedieme odhady regresných koeficientov získané metódou najmenších štvorcov:

Regr. koeficient	MNČ
$\beta_0$	1.20793
$\beta_1$	1.99619
$\beta_2$	-0.04609

**Tabuľka 3.** Výsledky odhadov regresných koeficientov MNČ pre TP

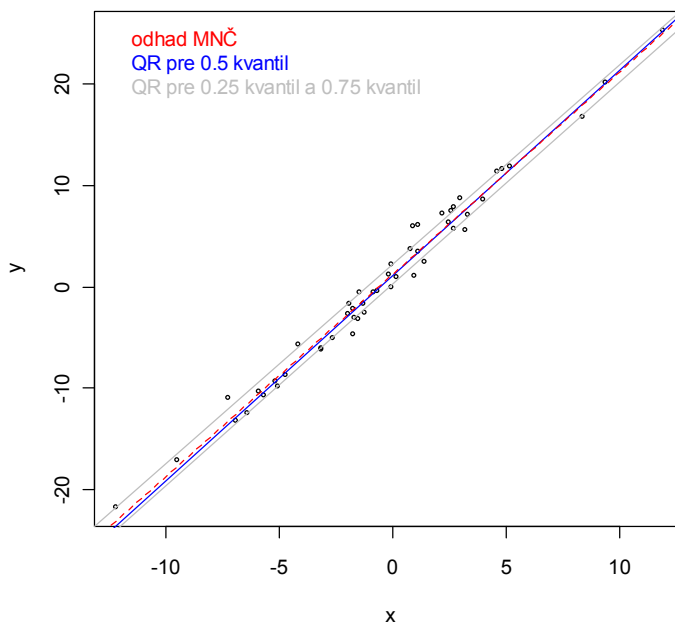
Zmeny regresných koeficientov v závislosti na jednotlivých kvantiloch si znázorníme graficky:



**Obrázok 6.** Odhady regresných koeficientov pre TP

Červené konštantné priamky predstavujú odhad koeficientov metódou najmenších štvorcov a 90% interval spoľahlivosti. Šedá oblasť predstavuje 90 % intervaly spoľahlivosti pre jednotlivé regresné koeficienty v závislosti od kvantilov.

Následne si vykreslíme odhady pre regresné kvantily a porovnáme ich s odhadom získaným metódou najmenších štvorcov.



**Obrázok 7.** Kvantilové priamky pre TP

Teraz predvedieme niektoré aplikácie regresných kvantilov pri testovaní hypotéz. Prirodzenou otázkou je, či odhadnuté kvantilovo regresné vzťahy potvrdzujú tzv. hypotézu o posunutí polohy (location shift hypothesis), ktorá predpokladá, že všetky kvantilové funkcie majú rovnaké stĺpcové parametre. Jedná sa o test zhody smerníc a teda testujeme platnosť

$$H_0: \beta_i(\tau) = \alpha_i \text{ pre } i = 1, 2, \dots, n.$$

Dôsledkom platnosti tejto hypotézy by bola lineárna závislosť premennej  $Y$  na  $X$ .

Ukážeme si ako si overiť platnosť tejto hypotézy v kvantiloch. Prevedieme tri procesy odhadovania (fits):

```
fit1 <- rq(y ~ x, tau = 0.25),
fit2 <- rq(y ~ x, tau = 0.5),
fit3 <- rq(y ~ x, tau = 0.75),
```

na ktoré použijeme funkciu anova (viď [10] help pre funkciu anova.rq) :

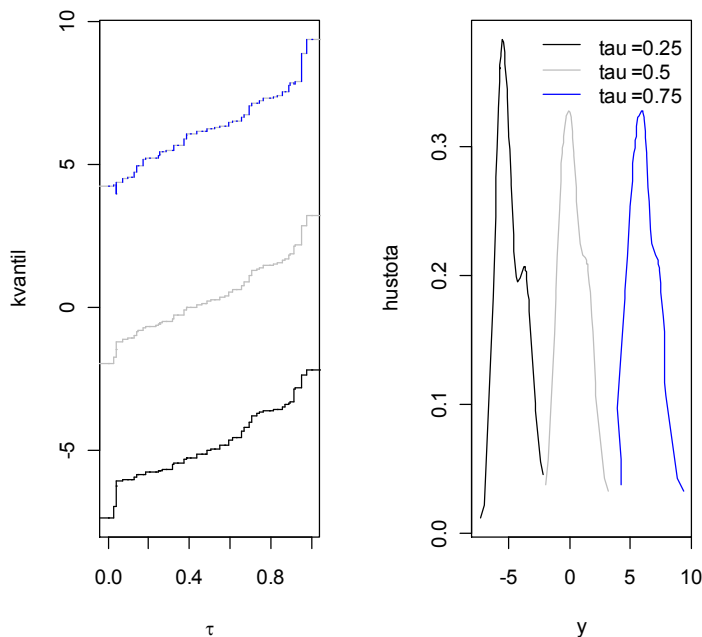
```
anova(fit1, fit2, fit3),
```

ktorá nám vráti tabuľku analýzy rozptylov :

```
Quantile Regression Analysis of Deviance Table
Model: y ~ x
Joint Test of Equality of Slopes: tau in { 0.25 0.5 0.75 }
  Df Resid Df F value Pr(>F)
1 2      148 0.6272 0.5355
```

Na základe výsledkov tohto testu nemôžeme zamietnuť platnosť hypotézy  $H_0$  a teda ani hypotézu o lineárnej závislosti skúmanej premennej  $Y$  na  $X$ .

Pre lepšie grafické pochopenie vykreslíme:



**Obrázok 8.** Odhadnuté podmienené kvantilové a hustotové funkcie pre  $Y$

Tento test nám neumožnil zamietnuť hypotézu  $H_0$  a preto skúsime previesť dôslednejší test, ktorý je založený na regresných poradiach a berie do úvahy hodnoty vo viacerých kvantiloch rozdelenia  $Y$  (nielen v kvartiloch):

```
KhmaladzeTest(y ~ x, taus = -1, nullH = "location")
```

výsledkom ktorého je:

```
$nullH
[1] "location"

$Tn
[1] 1.035798
```

Po dohľadani v tabuľke obsahujúcej kritické hodnoty tohto testu dostávame, že 1% kritická hodnota pre naše dáta ( $p = 2$ ,  $\epsilon = 0.05$ ) je rovná 4,199 (viď [1] tab. B2) a teda ani týmto testom nezamietame hypotézu o zmene polohy a tak nemôžeme vylúčiť, že  $Y$  je lineárne závislé na  $X$  a teda nami navrhovaná voľba modelu

$$Y = \beta_0 + \beta_1 \cdot X$$

sa javí ako vhodná.

## 4.2. Aplikácia prevedená na reálnych dátach

Praktické použitie metód kvantilovej regresie si ukážeme na príklade, v ktorom budeme na reálnych dátach zozbieraných v rokoch 1993-94 z 1283 amerických univerzít, určovať závislosť výšky školného príslušnej univerzity, na premenných uvedených v tabuľke 3. Túto závislosť sa pokúsime slovne interpretovať a zdôvodniť vývoj regresných koeficientov jednotlivých premenných v závislosti na kvantile  $\tau$ . Taktiež overíme, či bolo vhodné voliť metódy lineárnej regresnej analýzy a zameriame sa na to ako sa mení závislosť školného na svojich regresoroch v závislosti od toho, či je škola štátna, alebo súkromná .

Datový súbor bol získaný z dvoch zdrojov a to :

- 1) U.S. News & World Report's 1995 Guide to Americas Best Colleges
- 2) AAUP's (American Association of University Professors) 1994 Salary Survey

Názov premennej	stručná charakteristika
<i>tuition</i>	výška vysokoškolského školného v \$
<i>pcttop25</i>	percento novoprijatých študentov patriacich k 25% najlepších zo SŠ
<i>sf_ratio</i>	pomer členov fakulty k študentom
<i>fac_comp</i>	priemerná kompenzácia školného
<i>accrate</i>	podiel prijatých študentov k prihláseným
<i>graduat</i>	percento študentov, ktorí úspešne doštudujú
<i>pct_phd</i>	percento vyučujúcich s Ph.D
<i>fulltime</i>	percento študujúcich na dennom štúdiu
<i>alumni</i>	percento absolventov univerzity, ktorí jej prispievajú
<i>num_enrl</i>	počet novoprijatých študentov
<i>public/private</i>	premenná určujúca, či je škola štátna , alebo súkromná (štátna=0 ,súkromná =1)

**Tabuľka 4.** Charakteristiky univerzít

Uvažujeme regresný model v tvare

$$\begin{aligned}
 & \textit{tuition} = \\
 & \beta_0 + \beta_1 \textit{pcttop25} + \beta_2 \textit{sf\_ratio} + \beta_3 \textit{fac\_comp} + \beta_4 \textit{accrate} + \beta_5 \textit{graduat} \quad (14) \\
 & + \beta_6 \textit{pct\_phd} + \beta_7 \textit{fulltime} + \beta_8 \textit{alumni} + \beta_9 \textit{num\_enrl} + \beta_{10} \textit{public} / \textit{private}.
 \end{aligned}$$

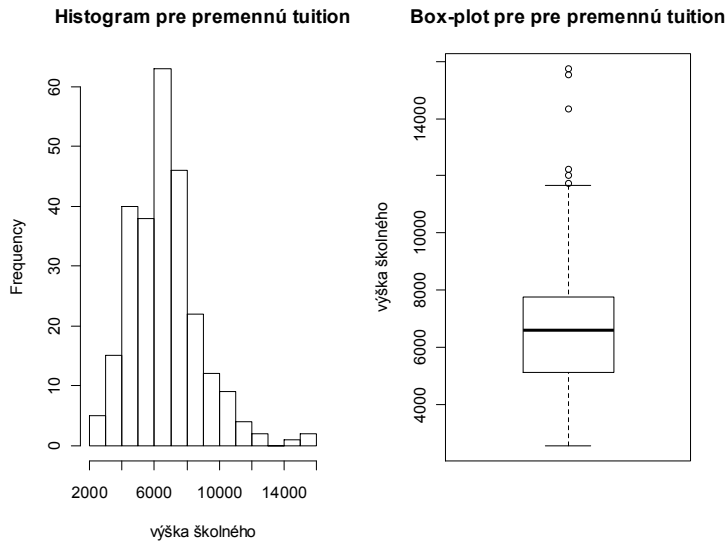
Náš datový súbor obsahuje informácie z 1283 amerických univerzít, avšak nie ku každej univerzite sú uvedené kompletne údaje . Použitím príkazu

```
data<- dataset[complete.cases(dataset), ]
```

vyradíme zo súboru univerzity, ku ktorým neboli dodané kompletne údaje pre všetky premenné. Takto vzniklý datový súbor obsahuje 804 univerzít s údajmi pre všetky premenné.

Bližší náhľad na hodnoty hodnoty výšky školného *tuition* :

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
2340	7146	9842	10290	12850	25750



**Obrázok 9.** Histogram a boxplot hodnôt *tuition*

Z boxplotu vidíme, že náhodný výber obsahuje viaceré odľahlé pozorovania .

Ešte skôr, než posúdime odhady regresných koeficientov  $\beta_i, i = 1, \dots, 10$ , otestujeme jednotlivé regresory na nevýznamnosť. K tomuto účelu vyberieme jednu z metód, ktoré sme si predviedli v teoretickom príklade. Ako najvhodnejší sa pri väčšom počte regresorov javí postup pomocou stepwise procedúry, ktorej priebeh sme si popísali v 4.1. Na začiatku uvažujme náš pôvodný regresný model (14) pomocou stepwise procedúry sa budeme snažiť vybrať jeho pre regresnú analýzu čo najvhodnejší submodel.

Výpočet v programe R:

```
max.model<-as.formula(paste("tuition~",
paste(names(+ dataset[c(2:11)]),collapse="+"))
model<-step(rq(tuition~1,method="fn"),scope=max.model,trace=0)
```

Navrhnutý submodel má tvar:

$$\text{tuition} = \beta_0 + \beta_1 \text{sf\_ratio} + \beta_2 \text{fac\_comp} + \beta_3 \text{graduat} + \beta_4 \text{pct\_phd} + \beta_5 \text{fulltime} \\ + \beta_6 \text{alumni} + \beta_7 \text{public / private}$$

Na tento model aplikujeme príkaz:

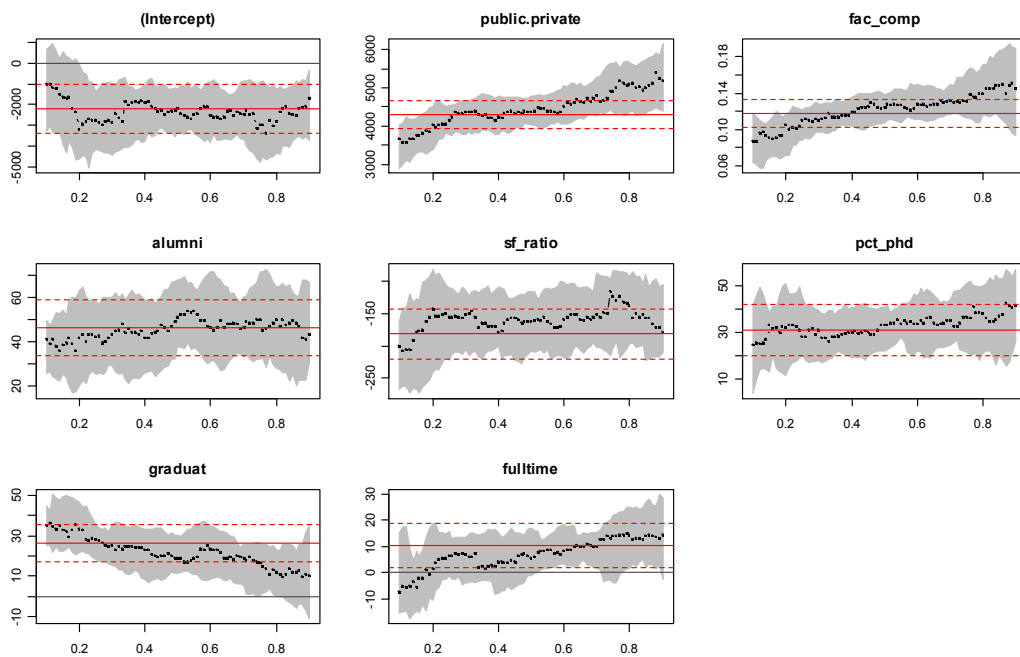
```
regresory <- summary(rq(model, tau = c(0.25, 0.5, 0.75)))
```

Výsledky odhadov regresných koeficientov pre najvýznamnejšie kvantily spolu s 95 % intervalmi spoľahlivosti si uvedieme v tabuľke:

<b>Regresný koeficient</b>	<b>0.25</b>	<b>0.5</b>	<b>0.75</b>
$\beta_0$	-2774.26811 (-4149.52193, -1351.57218)	-2221.79773 (-3337.04237, -1239.51208)	-2976.22398 (-4756.54820, -905.21631)
$\beta_1$	-155.36114 (-223.05887, -111.13941)	-163.30405 (-208.00637, -117.28089)	-123.47423 (-204.37742, -84.28445)
$\beta_2$	0.10948 (0.08882, 0.12202)	0.12382 (0.10803, 0.13873)	0.13053 (0.11300, 0.16339)
$\beta_3$	27.91620 (18.34451, 35.18609)	18.60247 (11.29319, 30.04818)	16.67539 (0.25042, 26.30793)
$\beta_4$	27.93436 (19.98321, 45.64387)	33.05932 (22.14588, 40.67896)	36.45428 (21.41629, 51.06758)
$\beta_5$	6.28444 (-1.42940, 15.53886)	5.44838 (-1.42445, 16.97730)	13.71562 (-3.87215, 21.62930)
$\beta_6$	42.17818 (25.28324, 58.31736)	51.31407 (37.96737, 68.00181)	50.12455 (31.98825, 71.62279)
$\beta_7$	4145.92048 (3785.38198, 4546.09932)	4359.04528 (4033.89952, 4773.58586)	4909.94529 (4371.58425, 5629.10246)

**Tabuľka 5.** Výsledky odhadov regresných koeficientov pre PP

To ako sa menia regresné koeficienty v závislosti na jednotlivých kvantiloch pre výšku školného si znázorníme graficky :



**Obrázok 10.** Odhady regresných koeficientov pre príklad s reálnymi dátami

Teraz sa pokúsime slovne interpretovať výsledky kvantilovej regresnej analýzy založené na grafoch obsiahnutých v obrázku (14). Pripomeňme, že niektoré premenné boli pri vyberaní najvhodnejšieho submodelu z regresnej analýzy vyradené a tak ich regresné koeficienty uvažujeme rovné nule pre všetky kvantily.

- **Regresor : public.private , regresný koeficient :  $\beta_7$**   
premenná určujúca, či je škola štátna, alebo súkromná štátna=0, súkromná=1

Školné v súkromných školách je, už z výsledkov získaných metódou najmenších štvorcov, očividne vyššie než v školách štátnych a to v priemere o 4359 dolárov. Regresná analýza nám však navyše poodhaľuje, že tento rozdiel je oveľa nižší pre nízke kvantily a podstatne vyšší na opačnom chvoste rozdelenia.

- **Regresor : fac\_comp , regresný koeficient :  $\beta_2$**   
priemerná kompenzácia školného

Regresor sa javí na základe porovnaní regresných koeficientov ako málo významný, musíme však vziať do úvahy, že kým hodnoty ostatných regresorov predstavujú percentá prípadne nadobúdajú hodnoty  $\{0,1\}$ , tak tento regresor nadobúda rádovo vyššie hodnoty a tak je jeho vplyv aj pri nízkom regresnom koeficiente na výšku školného značný. Tento vplyv sa s narastajúcim školným postupne zvyšuje.

- **Regresor : alumni , regresný koeficient :  $\beta_6$**   
percento absolventov univerzity, ktorí je prispievajú

Vplyv tohto regresora nám osciluje okolo priemeru získaného metódou najmenších štvorcov.

- **Regresor : sf\_ratio , regresný koeficient :  $\beta_1$**   
pomer členov fakulty k študentom

Pomer študentov k vyučujúcim má výrazný negatívny dopad na výšku školného. Tento negatívny dopad sa trochu prekvapujúco s rastúcim kvantilom školného znižuje, čo si ale môžeme vysvetliť pomocou regresora udávajúceho percento vyučujúcich s Ph.D., ktorý v závislosti na školnom stúpa a tak môže aj menší počet členov fakulty zabezpečiť kvalitné vzdelanie. Pre vysoké kvantily sa tento trend obracia a negatívny dopad na školné sa zvyšuje.

- **Regresor : pct\_phd , regresný koeficient :  $\beta_4$**   
percento vyučujúcich s Ph.D

Vplyv tohto regresora je kladný a postupne s rastúcim školným pomerne významne rastie. To môže byť spôsobené tým, že u škôl s vysokým pct\_phd je nízke sf\_ratio a výdaje na členov fakulty s Ph.D predstavujú vysoké náklady pre univerzitu, ktoré musia byť kompenzované výškou školného.

- **Regresor : *graduat* , regresný koeficient :  $\beta_3$**   
percento študentov, ktorí úspešne doštudujú

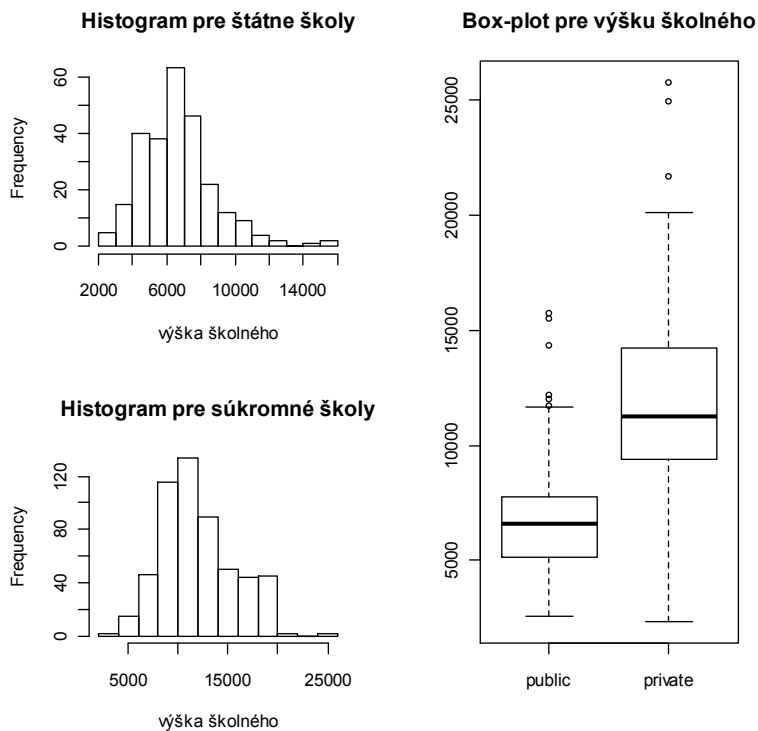
Vplyv tohto regresoru postupne s výškou školného klesá, čo si môžeme vysvetliť tým, že školy poskytujúce kvalitnejšie vzdelanie, a tým pádom majúce vyššie školné, sú pre študentov náročnejšie, čo znižuje podiel študentov, ktorí školu úspešne doštudujú. Avšak to týchto študentov natoľko neovplyvňuje pretože táto náročnosť zvyšuje prestíž školy a takto získaného titulu a teda pozitívny vplyv podielu doštudovaných študentov je zmenšený úbytkom prestíže.

- **Regresor : *fulltime* , regresný koeficient :  $\beta_5$**   
percento študujúcich na dennom štúdiu

Vplyv tohto regresoru nám postupne s výškou školného rastie, a kým u škôl z nízkych kvantilov môže byť tento vplyv dokonca záporný, pretože tieto školy s predpokladanou nižšou kvalitou sú zamerané na externých študentov, tak u vyšších kvantilov je vplyv podielu denných študentov na výšku školného už kladný.

Pri skúmaní výšky školného sa vynára otázka, či a o koľko sa líši školné v súkromných a štátnych školách a či sa tento rozdiel odráža aj v hodnotách charakterizujúcich premenných. Naším cieľom bude určiť ako sa líši vplyv charakterizujúcich premenných v jednotlivých kvantiloch v závislosti od toho, či sa jedná o školu štátnu alebo súkromnú.

Najprv si graficky znázorníme hodnoty výšky školného rozdelené podľa toho, či sa jedná o štátne, alebo súkromné školy:



**Obrázok 11.** Histogramy a boxploty pre hodnoty *tuition* rozdelené podľa premennej *public/private*

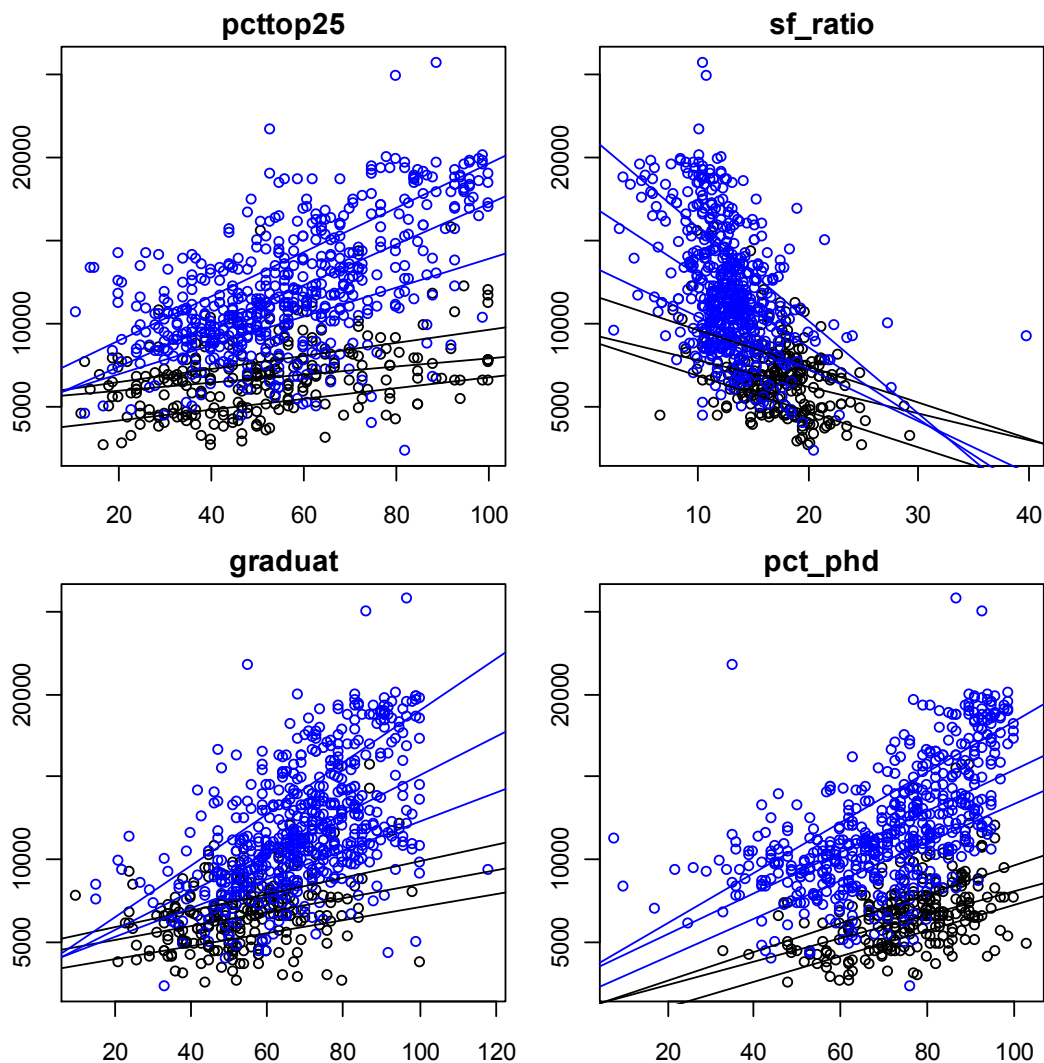
Z grafov vidíme, že školné v súkromných školách:

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
2340	9400	11280	12000	14200	25750

je zvyčajne vyššie ako v školách štátnych:

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
2580	5130	6594	6698	7765	15730

Teraz si vykreslíme grafy závislosti premennej *tuition* na vybraných charakterizujúcich premenných (*pcttop25*, *sf\_ratio*, *graduat*, *pct\_phd*), od ktorých vo všeobecnosti očakávame, že popisujú kvalitu vzdelania poskytnutého vysokou školou. Data rozdelíme do dvoch skupín podľa toho, či sa jedná o školu štátnu (čierna farba), alebo školu súkromnú (modrá farba) a do grafov zakreslíme regresné kvantily v kvartiloch pre príslušné skupiny.



**Obrázok 12.** Vykreslené závislosti *tuition* na vybraných regresoroch

Z grafov jednoznačne vyplýva, že v súkromných školách majú tieto regresory na stanovenie výšky školného väčší efekt ako v školách štátnych.

Na záver prevedieme opäť test hypotézy o zmene polohy (location shift), ktorou si overíme, či bolo vhodné modelovať výšku školného lineárnym regresným modelom.

Budeme overovať pre jednotlivé regresory ( $i = 1, 2 \dots, 7$ ) hypotézu

$$H_0: \beta_i(\tau) = \alpha_i$$

a následne spoločnú hypotézu pre všetky regresory:

$$H_0: \beta_i(\tau) = \alpha_i \text{ pre } i = 1, 2 \dots, 7.$$

To dosiahneme použitím príkazu v  $\mathbb{R}$

```
KhmaladzeTest(model, taus = -1, nullH = "location").
```

Výsledok testu si uvedieme v tabuľke:

<i>Premenná</i>	<i>T</i>
public.private	1.6879474
fac_comp	2.0102988
alumni	0.4734559
sf_ratio	1.0721061
pct_phd	0.5140671
graduat	0.5944883
fulltime	0.7793660
<b>spoločný efekt</b>	<b>6.523212</b>

**Tabuľka 6.** Testy hypotézy o zmene polohy

Po dohľadani v tabuľke obsahujúcej kritické hodnoty tohto testu dostávame, že 5% kritická hodnota pre spoločný efekt ( $p = 8$ ,  $\epsilon = 0.05$ ) je rovná 9,552 (viď [1] tab. B2) a teda nezamietame hypotézu o zmene polohy a teda nemôžeme vylúčiť, že  $Y$  je lineárne závislé na vektore  $\mathbb{X}$  a teda voľba lineárnej regresie sa javí ako vhodná a ani štatistiky pre jednotlivé kovariáty nie sú zamietajúce na 5% hladine.

## Záver

V tejto bakalárskej práci sú obsiahnuté teoretické základy, z ktorých bola odvodená teória regresných kvantilov. Snažili sme sa o ucelenejší pohľad na teóriu kvantilovej regresie, z pohľadu začlenenie odhadov v nej vzniklých do rámca M-odhadov a ukázali sme jej súvis s metódami klasickej regresie, špeciálne s odhadmi prostredníctvom metódy najmenších štvorcov. Poukázali sme na jeden z prístupov v lineárnom programovaní pre výpočet regresných koeficientov a predviedli sme odvodenie jeho algoritmu (Frisch-Newtonov algoritmus). Ukázali sme si štatistické použitie regresných poradií, vzniklých pri výpočte duálnej úlohy k úlohe pre výpočet regresných koeficientov, a následne sme ho demonštrovali v praktickej časti. V tej sme ilustrovali znalosti nadobudnuté v teoretickej časti na dvoch príkladoch, ktoré sme volili takým spôsobom, aby z nich bolo jasne viditeľné, aký potenciál má použitie metód kvantilovej regresie a nakoľko nám skúmanie priebehu regresných koeficientov pre jednotlivé kvantily môže ozrejmiť závislosť medzi regresormi a skúmanou premennou. Teória kvantilovej regresie je oveľa obširnejšia a prináša viac aplikácii ako je uvedené v tejto práci, účelom ktorej bolo priniesť čitateľovi základné znalosti potrebné pre ďalšie štúdium a prehĺbiť jeho záujem o ich ďalšie štúdium.

## Zoznam použitej literatúry

- [1] Koenker R.: Quantile Regression. Cambridge University Press, New York, 2005.
- [2] Gutenbrunner C., Jurečková J.: Regression Rank Scores and Regression Quantiles. The Annals of Statistics, Volume 20, pp. 305–330, 1992.
- [3] Jurečková J.: Robustní statistické metody, Praha, 2001.
- [4] Blatná D.: Robustní přístup v lineární regresi, Vysoká škola ekonomická v Praze, Praha, 2008 pp. 255–265.
- [5] Hekimoglu S., Erenoglu R.C., Kalina J.: Outlier detection by means of robust regression estimators for engineering science. Journal of Zhejiang University – Science A (JZUS-A), 10 (6), pp. 909-921, 2009.
- [6] Saleh A.K.Md.E., Picek J., Kalina J.: Nonparametric estimation of regression parameters in measurement error models. Metron 67 (2), pp. 177-200, 2009.
- [7] Zvára K., Štěpán J.: Pravděpodobnost a matematická statistika, Matfyzpress, Praha, 2006.
- [8] Remo A.: Regresní kvantily, Diplomová práce, Masarykova univerzita v Brně, Přírodovědecká fakulta, Brno, 2008.
- [9] Jurík T.: Vybrané algoritmy lineárneho programovania, Písomná práca k dizertačnej skúške, Univerzita Komenského v Bratislave, Fakulta Matematiky, Fyziky a Informatiky, Bratislava, 2007.
- [10] Koenker R.: quantreg: Quantile Regression. R package version 4.67.  
<http://cran.r-project.org/package=quantreg>